

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problem Mailbox.**

THIS PAGE BLANK (USPTO)

PATENT ABSTRACTS OF JAPAN

AA

(11)Publication number : 10-247172

(43)Date of publication of application : 14.09.1998

(51)Int.Cl.

G06F 13/36

(21)Application number : 10-031923

(71)Applicant : COMPAQ COMPUTER CORP

(22)Date of filing : 05.01.1998

(72)Inventor : GOODRUM ALAN L

(30)Priority

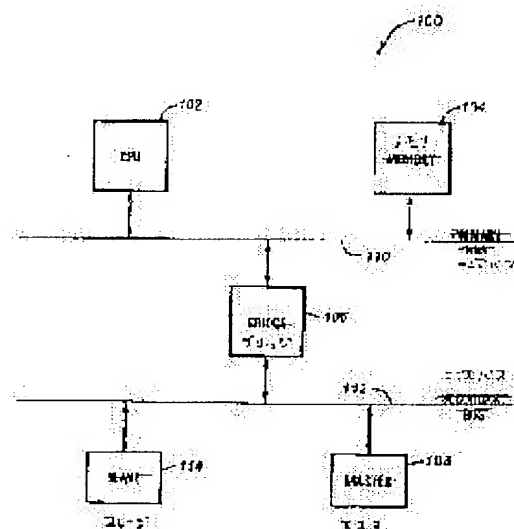
Priority number : 96 774746 Priority date : 31.12.1996 Priority country : US

(54) BUFFER RESERVATION METHOD FOR BUS BRIDGE SYSTEM

(57)Abstract:

PROBLEM TO BE SOLVED: To attain the optimum transfer of data between buses by receiving a reservation command from a bus master and reserving a data buffer for the bus master in response to the reservation command.

SOLUTION: A bridge 106 includes a mechanism to reserve a data buffer. Thus, a bus master 108 connected to the bridge 106 or an agent requests the bridge 106 to reserve a buffer for the data read transfer that is intended by the master 108. Then the bridge 106 gives a common interface to the device such as the master 108 existing on a secondary bus 112 and attains the communication with the device such as a memory 104 existing on a primary bus 110. The master 108 also can communicate with other secondary bus devices such as a bus slave 114, etc., without requesting any band at all against the bus 110.



LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

THIS PAGE BLANK (USPTO)

THIS PAGE BLANK (USPTO)

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平10-247172

(43) 公開日 平成10年(1998) 9月14日

(51) Int.Cl.⁹
G 0 6 F 13/36

識別記号
3 1 0

F I
G 0 6 F 13/36

3 1 0 F

審査請求 未請求 請求項の数36 O L 外国語出願 (全 76 頁)

(21) 出願番号 特願平10-31923

(22) 出願日 平成10年(1998) 1月5日

(31) 優先権主張番号 7 7 4 7 4 6

(32) 優先日 1996年12月31日

(33) 優先権主張国 米国 (U S)

(71) 出願人 591030868

コンパック・コンピューター・コーポレーション

COMPAQ COMPUTER CORPORATION

アメリカ合衆国テキサス州77070, ヒューストン, ステイト・ハイウェイ 249, 20555

(72) 発明者 アラン・エル・グッドラム

アメリカ合衆国テキサス州77375, トムボール, エイヴァンフィールド 16522

(74) 代理人 弁理士 社本 一夫 (外6名)

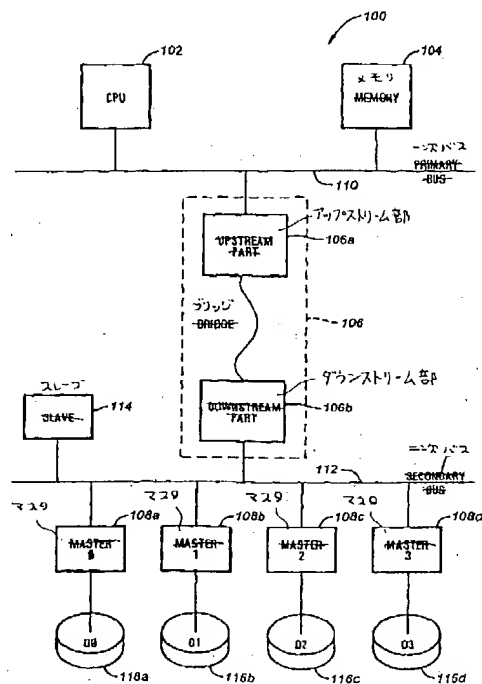
最終頁に続く

(54) 【発明の名称】 バス・ブリッジ・システムのためのバッファ・リザベーション方法

(57) 【要約】

【課題】 バス・ブリッジを通じたデータ転送方法を提供する。

【解決手段】 バス・ブリッジは、データを格納し、データを先取りし、データのライト・ポスティングを行うための多数のデータ・バッファを含む。バス・ブリッジと通信するデバイスは、2つの予約機構の1つによってバッファを予約することができる。予約機構は、バス・ブリッジに、アドレスおよびバイト・カウントを供給する。また、予約はいずれの上流側バス・ブリッジにも送出することができる。予約されたバッファは、バス・アクセスの効率的な使用のために先取りされる。バッファが予約されていない場合、データの先取りおよびフラッシュは、別のアルゴリズムにしたがって行われる。



【特許請求の範囲】

【請求項1】 バス・ブリッジとバス・マスタの間における通信方法であって、前記バス・ブリッジは少なくとも1つのデータ・バッファを有する、方法において、

(a) 前記バス・マスタから予約コマンドを受け取るステップと、

(b) 前記予約コマンドにตอบสนองして、前記バス・マスタのためにデータ・バッファを予約するステップと、を備える方法。

【請求項2】 請求項1記載の方法において、(c) ステップ(a)の後に、バス・コマンドを受け取るステップを更に備える方法。

【請求項3】 請求項2記載の方法において、前記バス・コマンドは開始アドレスを含み、前記バス・コマンドおよび前記予約コマンドは共になってアドレス範囲を定義する、方法。

【請求項4】 請求項3記載の方法において、前記予約コマンドはカウントを含む、方法。

【請求項5】 請求項3記載の方法において、前記バス・コマンドがリード動作である場合に、

(d) 前記リード動作に従ってデータを取り込むステップであって、該データは予約された前記バッファを通じて受け渡される、ステップと、

(e) 前記リード動作および前記アドレス範囲に従ってデータを先取りするステップであって、先取りした該データは前記の予約されたデータ・バッファに書き込まれる、ステップと、

(f) 前記の取り込んだデータおよび先取りしたデータを前記バス・マスタに供給するステップと、を更に備える方法。

【請求項6】 請求項5記載の方法において、

(g) 前記アドレス範囲を使い切るまで、ステップ

(e) および(f)を繰り返すステップと、

(h) ステップ(g)の後にデータ・バッファ予約を取り消すステップと、を更に備える方法。

【請求項7】 請求項1記載の方法において、前記予約コマンドは、フレーム信号がアサートされた後に受け取られる、方法。

【請求項8】 請求項1記載の方法において、前記バス・ブリッジおよび前記バス・マスタは、周辺素子相互接続(PCI)バスに結合するためのものである、方法。

【請求項9】 第1と第2バス・ブリッジの間における通信の方法であって、前記第2バス・ブリッジは前記第1バス・ブリッジとバス・マスタとの間を結合するものである、方法において、

(a) 前記第2バス・ブリッジにおいて前記バス・マスタから予約コマンドを受け取るステップと、

(b) 前記予約コマンドを前記第1バス・ブリッジに送出するステップと、を備える方法。

【請求項10】 請求項9記載の方法において、前記第2バス・ブリッジは、予約コマンドが送出されるか否かを選択するためのプログラム可能な送出ビットを含み、更に、

(c) ステップ(b)の前に予約コマンドの送出をイネーブルするために前記送出ビットがセットされているか否かについて判定を行うステップを含み、

ステップ(b)は、送出をイネーブルとするために前記送出ビットがセットされることを条件とする、方法。

【請求項11】 請求項9記載の方法において、

(d) ステップ(a)の後にバス・コマンドを受け取るステップと、

(e) 前記予約コマンドおよび前記バス・コマンドにตอบสนองして、前記バス・マスタのためにデータ・バッファを予約するステップと、を更に備える方法。

【請求項12】 請求項11記載の方法において、前記バス・コマンドは開始アドレスを含み、前記バス・コマンドおよび前記予約コマンドは共になってアドレス範囲を定義する、方法。

【請求項13】 請求項12記載の方法において、前記予約コマンドはカウントを含む、方法。

【請求項14】 請求項12記載の方法において、前記バス・コマンドがリード動作である場合に、

(f) 前記リード動作に従ってデータを取り込むステップであって、該データは予約された前記バッファを通じて渡される、ステップと、

(g) 前記リード動作および前記アドレス範囲に従ってデータを先取りするステップであって、先取りした該データは前記の予約したデータ・バッファに書き込まれる、ステップと、

(h) 前記の取り込んだデータおよび先取りしたデータを前記バス・マスタに供給するステップと、を更に備える方法。

【請求項15】 請求項14記載の方法において、

(i) 前記アドレス範囲を使い切るまでステップ(g)および(h)を繰り返すステップと、

(j) ステップ(i)の後に前記のデータ・バッファ予約を取り消すステップと、を更に備える方法。

【請求項16】 請求項9記載の方法において、前記予約コマンドは、フレーム信号をアサートした後に受け取られる、方法。

【請求項17】 請求項9記載の方法において、前記バス・ブリッジおよび前記バス・マスタは、周辺素子相互接続(PCI)バスに結合するためのものである、方法。

【請求項18】 バス・ブリッジにおいてバッファを予約する方法であって、

(a) フレーム指示を与えるステップと、
(b) 前記フレーム指示をあたえるときに予約コマンドを与えるステップと、
(c) ステップ(b)の後でかつ前記フレーム指示を与えている間に、バス・コマンドを与えるステップと、
を備える方法。

【請求項19】 請求項18記載の方法において、前記バス・コマンドは開始アドレスを含み、前記バス・コマンドおよび前記予約コマンドは共にアドレス範囲を定義する、方法。

【請求項20】 請求項18記載の方法において、前記予約コマンドはカウントを含む、方法。

【請求項21】 コンピュータ・システムであって、主メモリと、
前記主メモリに結合されるプロセッサと、
大容量記憶システムと、
前記大容量記憶システムに結合されるバス・マスタであって、予約コマンドおよびバス・コマンドを含むバス動作を与えるように動作可能なバス・マスタと、
前記主メモリと前記バス・マスタの間に結合されるバス・ブリッジであって、少なくとも1つの割り当て可能なデータ・バッファを有し、前記バス動作を受け取り、該バス動作にตอบสนองして、前記バス・マスタの排他的使用のためにデータ・バッファを予約するように動作可能であるバス・ブリッジと、
を備えるコンピュータ・システム。

【請求項22】 請求項21記載のコンピュータ・システムにおいて、前記バス・コマンドは開始アドレスを含み、前記バス・コマンドおよび前記予約コマンドは共にアドレス範囲を定義する、コンピュータ・システム。

【請求項23】 請求項22記載のコンピュータ・システムにおいて、前記予約コマンドはカウントを含む、コンピュータ・システム。

【請求項24】 請求項22記載のコンピュータ・システムにおいて、前記バス・コマンドがリード動作である場合に、
前記バス・ブリッジは、前記リード動作に従ってデータを取り込み、前記アドレス範囲に従ってデータを先取りし、データは要求されたときに前記バス・マスタに供給される、コンピュータ・システム。

【請求項25】 請求項24記載のコンピュータ・システムにおいて、前記バス・マスタが前記アドレス範囲の最後のアドレスにおいて読み取りを行うときに、前記バス・ブリッジは前記データ・バッファの予約を取り消す、コンピュータ・システム。

【請求項26】 請求項21記載のコンピュータ・システムにおいて、前記予約コマンドは、フレーム信号がアサートされた後に受け取られる、コンピュータ・システム。

【請求項27】 請求項21記載のコンピュータ・システムにおいて、前記バス・ブリッジおよび前記バス・マスタは周辺素子相互接続(PCI)バスに結合される、コンピュータ・システム。

【請求項28】 請求項21記載のコンピュータ・システムにおいて、前記バス動作は、前記大容量記憶システムに書き込むためのメモリ・リード動作を含む、コンピュータ・システム。

【請求項29】 コンピュータ・システムであって、主メモリと、
前記主メモリに結合されるプロセッサと、
大容量記憶システムと、
前記大容量記憶システムに結合されるバス・マスタであって、予約コマンドとバス・コマンドとを含むバス動作を与えるように動作可能なバス・マスタと、
前記主メモリに結合される第1バス・ブリッジと、
前記第1バス・ブリッジと前記バス・マスタとの間に結合される第2バス・ブリッジであって、少なくとも1つの割り当て可能なデータ・バッファと送出ビットとを有し、前記バス動作を受け取り且つ前記バス動作にตอบสนองして前記バス・マスタの排他的使用のためにデータ・バッファを予約するように動作可能であり、前記送出ビットがセットされている場合に前記予約コマンドを前記第1ブリッジに送出するように動作可能である、第2バス・ブリッジと、
を備えるコンピュータ・システム。

【請求項30】 請求項29記載のコンピュータ・システムにおいて、前記バス・コマンドは開始アドレスを含み、前記バス・コマンドおよび前記予約コマンドは共にアドレス範囲を定義する、コンピュータ・システム。

【請求項31】 請求項30記載のコンピュータ・システムにおいて、前記予約コマンドはカウントを含む、コンピュータ・システム。

【請求項32】 請求項29記載のコンピュータ・システムにおいて、前記バス・コマンドがリード動作である場合に、
前記第2バス・ブリッジは、前記リード動作に従ってデータを取り込み且つ前記アドレス範囲に従ってデータを先取りし、データは要求されたときに前記バス・マスタに供給される、コンピュータ・システム。

【請求項33】 請求項32記載のコンピュータ・システムにおいて、前記バス・マスタが前記アドレス範囲の最後のアドレスにおいて読み取りを行うときに、前記第2バス・ブリッジは前記データ・バッファの予約を取り消す、コンピュータ・システム。

【請求項34】 請求項29記載のコンピュータ・システムにおいて、前記予約コマンドは、フレーム信号がアサートされた後に受け取られる、コンピュータ・システム。

【請求項35】 請求項29記載のコンピュータ・システムにおいて、前記第2バス・ブリッジおよび前記バス・マスタは周辺素子相互接続(PCI)バスに結合される、コンピュータ・システム。

【請求項36】 請求項29記載のコンピュータ・システムにおいて、前記バス動作は、前記大容量記憶システムに書き込むためのメモリ・リード動作を含む、コンピュータ・システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、一方のバスから他方のバスに動作を変換する方法に関し、更に特定すれば、バス間でデータを最適に転送するためのバッファ予約機構(buffer reservation mechanism)に関するものである。

【0002】

【従来の技術】パーソナル・コンピュータは、新たな技術が発展しコンピュータに組み込まれるに連れて、常に変化しつつある。マイクロプロセッサおよびメモリの性能向上の結果、コンピュータは非常に強力になり、以前では大型メインフレーム・コンピュータによってのみ処理可能であったタスクも、今や処理可能となっている。しかしながら、メインフレーム・コンピュータに完全に取って代わるには、コンピュータは心臓部のI/O(入力/出力)サブシステムが大量のメモリおよび格納容量に対応していなければならない。

【0003】ISA(Industry Standard Architecture:業界標準アーキテクチャ)、EISA(Extended Industry Standard Architecture:拡張業界標準アーキテクチャ)、およびPCI(Peripheral Component Interface:周辺素子インターフェース)を含むいくつかの標準化されたI/Oバスが、システム設計者には使用可能である。今日のコンピュータは、典型的に、これら3種類のいずれかを組み合わせることによって設計される。これらバス間でデータを移動するためには、ブリッジ・デバイス(bridge device)が典型的に設けられる。

【0004】ブリッジ・デバイスは、双方のバスを接続し、これらバス間でデータを転送すると共に、バス制御信号の変換を行う。バスは異なるものとすることができ、あるいは、特にPCIの場合には、ブリッジは単に同一の論理バスに対する電氣的な拡張(エクステンション)を与えることができる。この電氣的な分離のために、一方のバス・セグメント上における物理的デバイス数を制限しつつ、同時にPCIエージェントの総数は制限しない、PCIバスの要件を満たすことができる。PCIバスの用語では、エージェントとは、マスタ・デバイスおよびスレーブ・デバイスを含む、バスに接続する種類のデバイスを示す言葉である。

【0005】これらのバスは、全て、バス・マスタリング(bus mastering)と呼ばれる方式に対応しており、通

常プロセッサ以外のデバイスまたはエージェントが、アービタ(arbiter)にバスの制御を要求することができ、アービタがエージェントに制御を許可した場合、このエージェントはバス・マスタとなる。バス・マスタはその処理を、スレーブと呼ばれる別のエージェントに差し向ける。次いで、バス・マスタは、プロセッサの介入なしに処理を実行することができ、プロセッサを伴う場合よりも格段に効率的である。多くの場合、発生元バス(origination bus)上のマスタは、宛先バス(destination bus)上のスレーブと通信する。これらのバスがブリッジで結合されている場合、ブリッジの性能がマスタとスレーブ間の通信性能に大きく影響する。従って、この経路を最適化することが望ましい。

【0006】PCIバスの場合、性能を向上させる方法の1つは、遅延トランザクション又はリード・ポスティング(read posting)を可能にすることである。PCIバスおよび遅延トランザクションに関する更なる詳細は、オレゴン州ヒルズボロ(Hillsboro)のPCIスペシャル・インタレスト・グループ(Special Interest Group)から入手可能な“PCIローカル・バス仕様バージョン2.1(PCI LocalBus Specification version 2.1)”に記載されている。PCIローカル・バス仕様およびそれに関連する文書をここに参照として援用する。遅延トランザクションは、低速デバイスがリクエスト(要求)に応答してデータを用意している間、バスの使用を可能にする。従って、低速デバイスが待ち時間をバスに加える代わりに、バスを他の要求のために用いることができる。ブリッジにとって、宛先バスは、元の要求のターゲット(target)として動作しなかったインターフェースを意味する。遅延トランザクションは、マスタによる要求、ターゲットによる要求の完了、およびマスタによるトランザクションの完了という3つの段階(フェーズ)で、完了に至る。

【0007】第1フェーズでは、マスタはバス上にトランザクションを発生し、ターゲットがアクセスをデコードし、アクセスを完了するために必要な情報をラッチし、リトライ終了(retry-termination)によって要求を終了する。マスタは、遅延トランザクション終了を用いてトランザクションを完了しつつあるターゲットと、現時点において単にトランザクションを完了できないターゲットとの区別ができないので、要求を再度発行しなければならない。第2フェーズでは、ターゲットは、遅延要求からラッチした情報を用いて、宛先バス上で独立して要求を完了する。遅延要求が読み出し(リード(read))の場合、ターゲットは、要求データおよび完了ステータスを取得する。遅延要求が書き込み(ライト(write))の場合、ターゲットはライト・データを送り、完了ステータスを取得する。第3フェーズの間、マスタは首尾良くバスの再調停を行い、元の要求を再発行する。ターゲットはこの要求をデコードし、完了ステータス(お

よび、リード要求の場合はデータ)をマスタに供給する。発生元バス上で完了する前に宛先バス上で完了しなければならない全てのバス・コマンドは、遅延トランザクションとして完了させることができる。これらは、割込アクノレッジ(承認)、I/Oリード、I/Oライト、コンフィギュレーション・リード、コンフィギュレーション・ライト、メモリ・リード、メモリ・リード・ラインおよびメモリ・リード多重のコマンドを含む。メモリ・ライト、およびメモリ・ライト及び無効化(invalidate)のコマンドは、宛先バス上で完了する前に、発生元バス上で完了することができる。これらのコマンドは、遅延トランザクション終了を用いて完了されず、通常はポスト(post)される。

【0008】かかるブリッジ・デバイスの1つが、IntelのPCI-EISAブリッジ・チップ・セットである。82375EB/SB PCI-EISA Bridgeおよび82374EB/SB EISAシステムのコンポーネントは、対で動作してPCIバスを有するコンピュータにEISA I/Oインターフェースを提供する。チップ・セットは、PCIバス及びEISAバスの両方上のマスタ又はスレーブの何れでもよい。PCI-EISAデータ転送では、4つの32ビット・ポストッド・ライト・バッファ(32-bit posted write buffer)が設けられ、単一サイクルPCIバス・トランザクションの改善を図っている。EISA-PCIデータ転送では、4つの16バイト・ライン・バッファが含まれ、EISAバースト処理(EISA bursting)を支持する。

【0009】双方のバスを効率的に使用するために、殆どのブリッジは、当該ブリッジ自体に、ある量のデータ・バッファ機能を実施している。これによって、ブリッジはバス同士を互いに結合解除し、他方のバスによる速度低下を招くことなく、各バスをその最大速度で処理させることが可能となる。通常、ブリッジ内に実施可能なバッファには2種類があり、それらは、ライト・ポストイング・バッファ(write posting buffer)およびリード・プリフェッチ(先取り)・バッファ(read prefetch buffer)又はリード・アヘッド・バッファ(read ahead buffer)である。両種類共、いずれのバス上でも実施可能である。

【0010】ライト・ポストイング・バッファは一方のバスからのライト・データを受け入れ、当該バスに受け取りを知らせる。これによって、バスを解放し他のトランザクションを処理させる。ブリッジは、他のバスに書き込むことができるようになるまで、ライト・データを一時的に格納、又はポスト、する。リード先取りバッファは、単一のリード・アクセスからアドレスを取り込み、追加のデータも必要となることを推測してその読み取りを行う。次いで、ブリッジは、バッファが使用不可能になるかあるいはリード・アクセスによって使用されるまで、このデータをバッファ内に保持する。

【0011】インテル(Intel)のチップ・セットでは、バッファによって、ブリッジはピーク・データ転送速度で短いデータ・バーストを受け取ることができるようになる。例えば、EISAデバイスがPCIバス上のメモリからのデータを要求すると、ブリッジは、PCIメモリからそのバッファに4つの32ビット・データ・ワードをバースト状に入力し、EISAデバイスがバッファを読み取っている間、PCIバスを他のPCI要求元に解放することができる。したがって、PCIバスは、EISAデバイスにより遅延されない。バッファが満杯の間、EISAデバイスは、バッファからデータを読み取るように通知される。こうして、EISAデバイスは、そのバースト転送速度でバッファからデータを読み出すことができる。その後、ブリッジは、PCIバスを通じて短いデータ・バーストを実行しバッファを満杯に維持することによって、EISAバス転送速度を維持しようとする。しかしながら、他のEISAデバイスがメモリからのデータを要求した場合、バッファを空にして(フラッシュして)再度満たさなければならない。

【0012】また、バッファは、データの一貫性に伴う問題の原因にもなる。データをブリッジ内にバッファしている間、バス・エージェントおよびプロセッサが、メモリ内に実際に何が入っているかについて、異なる考えを有する場合がある。バス・マスタがブリッジを通じて、ブリッジの他方側のターゲットのメモリ範囲にリード要求を発行した場合、ブリッジは、処理能力およびデータ完全性(インテグリティ(integrity))という2つの矛盾する目標の間で、バランスを取らなければならない。処理能力のためにメモリ動作および転送を最適化する場合、ブリッジは、好ましくは、大きなブロックの先取りしたデータを保持するであろう。しかしながら、データを大きなブロックで先取りしたものの、使用されない場合には、ブリッジは、マスタによる後からのリード要求に対して、古いデータを供給することを防止しなければならない。加えて、使用されないデータの大きなブロックを先取りすると、バス帯域を浪費することにより、性能低下を招く。したがって、性能要件(performance concerns)を満たしつつ、データの完全性も保証する解決策を見いだすことが望ましい。

【0013】バスは異なる周波数でも動作可能である。したがって、ブリッジは、信号をある周波数の一方のバスから別の周波数の他方のバスに変換する際に、信号を同期させなければならない。連続データ・フローを達成すべき場合、同期は性能を犠牲にしなければならない。

【0014】

【課題を解決するための手段】2つのバスを接続するバス・ブリッジは、複数のデータ・バッファを含む。データ・バッファは、リード・データを格納するため、先取りしたデータ又はリード・アヘッド・データを格納するため、あるいはライト・ポストッド・データを格納する

ために用いられる。本発明によれば、これらデータ・バッファの使用によって、デバイスがバスにアクセスする回数が減少し、バスのレイテンシ(latency)が減少する。

【0015】ブリッジは、更に、データ・バッファをリザーブ(予約)する機構を含む。2つの代替機構が開示され、これにより、ブリッジに結合されているバス・マスタまたはエージェントは、ブリッジがバス・マスタの意図するデータ・リード転送のためにバッファを予約するように、要求することができる。ブリッジ階層がある場合、予約要求(リザーベーション・リクエスト)は、バッファ予約機構を支持するアップストリーム(上流)・ブリッジに渡される。

【0016】ブリッジが受け取る予約要求は、バス・マスタがデータを予約したデータ・バッファから除去する際に、ブリッジがデータを予め読み取る(リード・アヘッドする)又は先取り(プリフェッチ)するために必要なアドレッシング情報を含む。予約が行われる場合、3つの128バイト・データ・バッファがバス・マスタのために予約される。データは順番にそして巡回的に、マスタが読み取るに連れて、これらのデータ・バッファに先取りされる。データ・バッファが最初にマスタによってアクセスされるときに、次のデータ・バッファが先取りされる。アドレス範囲全体が読み取られ、データがバス・マスタによって除去されたときに、予約は取り消される。プログラム可能な時間期間内に、バッファがバス・マスタによってアクセスされない場合も、予約は取り消される。予約バッファを管理するための予約アルゴリズムを開示する。

【0017】ブリッジは、更に、予約機構がエージェントに使用可能でない場合に、バッファを管理する2つの代替アルゴリズムも含む。デフォルト・アルゴリズムを用いると、新しいバス・マスタが最初にバッファを予約せずにリードを行おうとした場合、バッファが使用可能であれば、ブリッジが3つの128バイト・データ・バッファから成るバッファをこのマスタに割り当てる。その他の場合、マスタは、リード動作をリトライするように強制される。一旦バッファが割り当てられたなら、ブリッジは要求されたデータを128バイト・データ・バッファの最初のものに読み込み、その後データをバス・マスタに供給する。マスタが、第1データ・バッファの終端を越えてバースト・リードを行い、マスタからの次のリード動作が、最後のリードが終了した位置から開始される場合、このマスタのために先取りがイネーブルされる。次のシーケンシャルのリードの時に、ブリッジは再び要求されたデータを第2データ・バッファに読み込み、その後それをマスタに供給する。この時点では、先取りは既にイネーブルされている。マスタが、バッファの1つの終端を越えて再度バースト・リードを行うときはいつでも、ブリッジは自動的に次の128バイト・バ

ッファに対してデータを先取りする。3つの128バイト・データ・バッファは、順番にそして巡回的に満たされる。バス・マスタが不連続なリードを要求した場合、またはCPUがマスタに影響を与え得るライト・トランザクションを実行した場合、バッファはフラッシュされる。

【0018】保存アルゴリズムは、デフォルト・アルゴリズムと同様であるが、ブリッジは先取りを行わない。バス・マスタがデータを全て読み取らない場合、マスタがバスからディスコネクト(接続解除)されるときに、残りのデータをフラッシュする。

【0019】

【発明の実施の形態】次に図1を参照すると、好適実施例によるブリッジ・デバイス106が図示されている。ブリッジ・デバイス106は、コンピュータ・システム100の一次バス110を二次バス112に結合する。コンピュータ・システム100は、更に、一次バス110に結合されているCPU102およびメモリ104を含む。好適実施例では、CPU102は、IBMコンパチブルのPCにおいて共通するような、インテルの486、ペンティアム7(Pentium7)または586級のプロセッサまたはそれらの同等物のようなプロセッサ(図示せず)、ならびに標準的なコンピュータ周辺装置および支援ロジック(図示せず)を含む。他のプロセッサや周辺素子も勿論使用可能であることは理解されよう。

【0020】ハード・ディスク・コントローラ、ビデオ・グラフィック・コントローラまたはネットワーク・インターフェース・コントローラのようなバス・マスタ108、およびシリアル(直列)・ポートまたはパラレル(並列)・ポートのようなバス・スレーブ114が、更に、二次バス112に接続されている。本発明において開示される原理の目的のためには、バス・マスタ108およびスレーブ114の具体的な機能性を示すことは不要であるので、ここでは、単に一般的なバス・マスタおよびスレーブとして言及するに止めておく。

【0021】一次バス110および二次バス112は、周辺素子相互接続(PCI)バスのような、標準的なI/Oバスである。尚、ここに開示する原理は、業界標準アーキテクチャ(ISA)、拡張業界標準アーキテクチャ(EISA)およびマイクロチャネル・アーキテクチャ(MCA: microchannel architecture)バスのような他のバスにも適用可能であることは理解されよう。更に、同じブリッジ設計原理は、PCIバスおよびEISAバスのような2種類の異なるバスを結合するために適用可能であることも、認めることができよう。好適実施例では、一次バスは、電氣的ファンアウト(electrical fanout)を限定したローカルPCIバスであり、ブリッジ106は、2次のPCIバス112によって一次バス110の機能性を拡張するために用いられる。

【0022】図2に示す図1の代替構成では、二次バス

112は、コンピュータ・システム100から離れた場所にある。更に、ブリッジ106は、アップストリーム（上流）部106aおよびダウンストリーム（下流）部106bを備える。ケーブル118がアップストリームとダウンストリームのブリッジ・インスタンスを互いに結合する。ダウンストリーム部106bは二次バス112に結合され、複数のバス・マスタ108aないし108dとの通信を行う。各バス・マスタは、コンピュータ・システム100に大容量記憶装置を提供するために、複数のハード・ディスク・ドライブ116aないし116dにも結合される。図1および図2の双方において、ブリッジ106は、バス・マスタ108のような二次バス上のデバイスに共通インターフェースを与え、メモリ104のような一次バス上のデバイスとの通信を可能にする。また、バス・マスタ108は、一次バス110からの帯域を全く要求することなく、バス・スレーブ114のような、その他の二次バス・デバイスとも通信することができる。したがって、変換ユニットおよびバス拡張部であることに加えて、ブリッジ106は分離バッファ(isolation buffer)としても機能する。

【0023】バス・マスタ108がメモリ104に対して読み取りまたは書き込みを行う場合のように、大きなデータ・ブロックをブリッジを通じて転送する際、ブリッジは、リード先取り(read prefetching)およびライト・ポスティング(write posting)のための多数のバッファを提供し、全体的なシステム性能向上を図る。これは、転送が発生し得る前に、双方のバスがブリッジ106と通信するという、典型的な要件をなくすのに役立つ。リード先取りの管理は、いくつかの代替的なアルゴリズムによって行われる。アルゴリズムの1つは、バス・マスタが非常に効率的なバス転送のためにリード先取りバッファを予約可能とする機能を含む。他の適応的アルゴリズムは、バス・マスタがバッファ予約を要求しない場合を扱う。しかしながら、これらの場合でも、リード先取りが高い成功確率を有することが所定の状態によって示された場合には、アルゴリズムはバッファを予約しデータを取引しようとする。

【0024】次に図3および図4に移ると、ブリッジ・デバイス106の第1および第2実施例のブロック図が示されている。図に見られるように、一次バスは一次バス・インターフェース200に結合し、二次バスは二次バス・インターフェース202に結合する。ブリッジ106は、入来データを受け取り、リード先取り又はリード・アヘッドを行い、また、ライト・ポスティングを行うためのバッファ領域204を含む。バッファは、処理の種類に応じたコンフィギュレーションに設定される。これについては、以下で説明する。

【0025】具体的には、ブリッジ106には、ライト動作の処理を向上させるために、4つの128バイト・ライト・ポスティング・バッファ218が設けられい

る。バス・マスタ108が主メモリ104に対してライト動作を実行する場合、データはライト・バッファ218に書き込み又はポスト(post)され、バス・マスタ108は、データが実際に主メモリ104に書き込まれる前に、ブリッジ106から完了の指示(indication)を受け取る。こうして、ブリッジ106は後から主メモリ104に対するライト動作を完了し、遅延を生じない。このように、実際の完了指示を待たないので、バス・マスタ108が二次バス102を制御する時間は大幅に短縮される。

【0026】また、バッファ領域204は、主メモリ104からのバス・マスタ・リード動作を向上するため、8つのデフォルト・バッファ210も含む。デフォルト・バッファ管理アルゴリズムがデフォルト・バッファ210を管理し、所定の条件の下において、リード・データの格納およびリード・データの先取りを行う。デフォルト・バッファ管理アルゴリズムは、できるだけ長く先取りデータを保持しながらも、古いデータを残す可能性を回避するという2つの間のバランスを取る。バッファ内の32ビット・ダブル・ワードのデータのチャンク(chunk)（以下DWORDSと呼ぶ）に対する3つの3ビット・ポインタを用いて、バッファを満杯にしたり空にしたりする。バッファ入力ポインタ(BUFF_IN_PNTR)は、次のDwordが格納される場所を指し示す。バッファ出力ポインタ(BUFF_OUT_PNTR)は、次のDwordが読み出される場所を指し示す。バッファ有効ポインタ(BUFF_VALID_PNTR)は、BUFF_IN_PNTRとBUFF_OUT_PNTRの間のどちらのアドレスが有効であることを示す。バスがロックされる状況を防止するために、ポストされたメモリ・ライト(posted memory write)の後に完了したリードは、バッファに格納することができるが、ポストされたメモリ・ライトが完了する後までは有効となることができないので、所定のアドレスのみが有効となるようにすればよい。

【0027】デフォルト・バッファ管理アルゴリズムの代わりとして、デフォルト・バッファ210を管理するために保存アルゴリズムが用意されている。保存アルゴリズムは先取りを行うことはなく、通常はバス・マスタによって要求されたデータを読み取るのみである。したがって、古いデータは完全に回避される。

【0028】更に、バッファ領域204は、主メモリ104からのバス・マスタ・リード動作の向上を図るために、8つの予約リソース（資源）220のプールも含む。これらの資源220は、以下で論ずる予約アルゴリズムによって管理される。各予約資源206は、3つの128バイト・データ・バッファと、対応する予約情報ブロック208で構成されている。各情報ブロック208は、予約バッファを制御するのに必要なレジスタを含む。かかるレジスタには、64ビットの予約の開始アドレスを収容する予約ベース・アドレス・レジスタ(RSRV_BASE_ADDR)、予約の最後のアドレスよりも1大きな64

ビット・アドレスを収容する予約制限アドレス・レジスタ(RSRV_LIMIT_ADDR)、非アクティブな予約資源を自動的に取り消すために4つの値(オフ、25ms、50ms、または100ms)の1つにプログラムされた破棄タイマ値、データが読み込まれた最後のライン・バッファへの5ビット・ポインタを収容する最後バッファ・インディケータ・レジスタ(LAST_BUFF)、各ライン・バッファの64ビット開始アドレスを収容するバッファ開始ポインタ・レジスタ(BUFF_START_PNTR)、ライン・バッファの最終アドレスよりも1大きな64ビット・アドレスを収容するバッファ制限ポインタ・レジスタ(BUFF_LIMIT_PNTR)、次のDwordが格納されるアドレスを保持するバッファ入力ポインタ・レジスタ(BUFF_IN_PNTR)、次のDwordが読み出されるアドレスを保持するバッファ出力ポインタ・レジスタ(BUFF_OUT_PNTR)、BUFF_IN_PNTRとBUFF_OUT_PNTRの間でどちらのアドレスが有効かを示すバッファ有効ポインタ・レジスタ(BUFF_VALID_PNTR)、バッファ状態レジスタ、およびバッファ先取りフラグ・レジスタが含まれる。

【0029】デフォルト・バッファ210および予約資源220は、各要求元マスタあるいは各プロセスまたはスレッドに割り当てられてもよく、こうすることにより、単一のマスタが多数のバッファを有することが可能となる。バッファ領域204はFIFOとして編成されており、勿論、最適なバス利用のために、マスタの要件にしたがってサイズを決定することができる。ブリッジ106は要求元およびそのターゲットには透明であることが好ましいので、バッファ領域204もトランザクションに対して透明であることが望ましい。しかしながら、バッファ領域204の動作を最適化するには、データが効率的に先取りされ、先取りされたデータが主メモリ内のデータとコヒーレントとなるようにバッファを制御することが望ましい。したがって、未処理のトランザクション(pending transaction)に対してバッファを最適化するには、マスタが読もうとしているアドレス範囲をブリッジが受け取ることが望ましい。こうすれば、一旦ブリッジがアドレス情報を取得したなら、ブリッジはアドレス範囲の最後まで自由にデータを先取りすることができる。予約資源を要求するための2つの代替案が用意されている。

【0030】図3に示す第1の代替的な予約機構において、二次インターフェース202はコマンド・デコード・ブロック212aを含む。メモリ・リード要求は下流側のバス・マスタからなので、コマンド・デコード・ブロック212aは二次インターフェース202に配置されているが、逆方向のリード動作では、コマンド・デコード・ブロックが一次インターフェース200において有効であることに留意されたい。予約資源は、主メモリ104から、二次バス112上に位置するバス・マスタから開始されるリード動作に対して最も効果的である

ので、以下の説明はこの文脈に沿って進めることにする。

【0031】予約資源220は、特殊なバス・コマンドがブリッジ106に受け取られた場合に、バス・マスタ108によって予約することができる。この点について、読者は図7も参照されたい。バス・マスタ108がターゲット・メモリ104からのリード・トランザクションのために予約資源220を予約したい場合、バス・マスタ108は、ターゲット・メモリ104から読み出される合計バイト数を示す予約コマンド(RSRV_COMMAND)を含む、第1リード・トランザクションを発行する。RSRV_COMMANDは、PCIバス上のコマンド/バイト・インエーブル(C/BE#)信号によって表される特殊なPCIバス・コマンドであり、図7のクロックの2に示すように、フレーム(Frame#)信号がアサートされトランザクションの開始が示されたときに、第1リード・トランザクション上でブリッジ106に提示される。また、予約コマンドは、AD[31:0]ライン上のカウントも含み、所望されるバイト数を示す。クロックの3以降では、PCIトランザクションは通常通りに進む。即ち、PCIトランザクションの種類が、アドレス/データ(AD)ライン上の開始アドレスと共に、C/BE#ラインに供給され、入手可能なときにデータが転送される(クロックの5、7および9に示す通り)。したがって、カウント(クロックの2)および開始アドレス(クロックの3)は、読み取るべきアドレス範囲を示す。このバッファを予約する範囲を示すには、開始アドレスおよび終了アドレスを示すことによる等のように、代替方法もあることは理解されよう。これらの代替方法は、予約コマンドのために考えられたものである。

【0032】図8に示す64ビットのアドレッシングでは、RSRV_COMMANDがデュアル(二重)・アドレス・コマンドおよび実際のリード・コマンドの前に位置する。二重アドレス・コマンドについての更なる詳細は、先に引用したPCI仕様に記載されている。

【0033】ブリッジ106がRSRV_COMMANDを受け取った場合、予約情報ブロック208のレジスタが初期設定される。開始アドレスはRSRV_BASE_ADDRに格納され、終了アドレスは、開始アドレスに合計バイト・カウントを加算することによって決定され、1だけ増分されてRSRV_LIMIT_ADDRに格納される。3つの128バイト・バッファの内1つが、開始アドレスをBUFF_START_PNTRに格納し、バッファ終了アドレスをBUFF_LIMIT_PNTRに格納することによって、初期設定される。また、3つのバッファ各々について、バッファ状態および先取りフラグも初期設定される。

【0034】一旦予約が行われたなら、バス・マスタ108はブリッジ106からディスクネクトし、後に再接続することにより、RSRV_COMMANDを再発行することなく処理を再開することができる。したがって、この代替案

では、単一のPCIバス・コマンドが、転送すべきバイト数を示す。一旦バッファが予約されたなら、データ・トランザクションは正常に進行するので、この追加情報に対する負担は1クロック・サイクルのみである。

【0035】ターゲットはRSRV_COMMANDを受け取らない。しかしながら、中間ブリッジが、下流(ダウンストリーム)側ブリッジからのRSRV_COMMANDを受け取る場合がある。図5および図6に示すように、ブリッジの階層構造がある場合、RSRV_COMMANDを支持する各ブリッジは、予約を上流(アップストリーム)に渡すことができる。各ブリッジのコンフィギュレーションでは、上流側ブリッジがバッファ予約を支持していない場合、コンフィギュレーション・ビット214(図3)がクリアされ、その他の場合、このビットはコンピュータの初期設定処理の間にコンフィギュレーション・ソフトウェアによってセットされ、予約(リザーベーション)は上流に渡される。図5は、ブリッジ階層の一例を示し、この場合、ブリッジ300は、バッファ予約に対応する(バッファ予約を支持する)2つのブリッジ106aと106bの間に挟まれた非予約ブリッジである。この例では、ブリッジ300がバッファ予約に対応していないので、ブリッジ106b内のコンフィギュレーション・ビットはクリアされ、RSRV_COMMANDは非予約ブリッジ300には送出されない。したがって、ブリッジ106aはバッファ予約に対応しているが、この特徴はブリッジ106aでは利用されない。

【0036】図6はブリッジ階層の一例を示し、この場合、最上位のブリッジ300はバッファ予約に対応していないが、2つの下位ブリッジ106aおよび106bはそれに対応している。この例では、ブリッジ300はバッファ予約に対応していない(バッファ予約を支持しない)ので、非予約ブリッジに最も近いブリッジ、即ち、ブリッジ106a内のコンフィギュレーション・ビットをクリアし、RSRV_COMMANDが非予約ブリッジ300に送出されないようにする。ブリッジ106b内のコンフィギュレーション・ビット214はセットされ、ブリッジ106bによって受け取られたRSRV_COMMANDがブリッジ106aに送出されるようにする。このようにして、非予約ブリッジは、RSRV_COMMANDが実施される場合に、バッファ予約の効果を切り離すように作用する。

【0037】次に、図4に示す第2の代替的な予約機構を参照すると、ブリッジ106は、特別なアドレスに対する処理をデコードするためのアドレス・デコード・ブロック212bを含む。コンピュータ・システム100のメモリ・アドレス・マップでは、この特別なメモリ・アドレスは、この予約機構のために割り当てられている。この特別アドレスに対するライト処理は、ブリッジ106による予約動作として解釈される。バス・マスタ108がリード動作のためにバッファ予約を望む場合、まず、2-3 32ビット(Dword)バースト・ライトを

特別アドレスに実行する。ブリッジ106はこのライト動作を横取り、コマンドを予約メールボックス(reservation mailbox)216に格納する。最初のDwordは、バス・マスタ108が読み取ろうとする合計バイト数である。2番目のDwordは開始アドレスの下位32ビットである。64ビット・アドレッシングを用いる場合、開始メモリ・アドレスの上位32ビットを収容する第3のDwordが書き込まれる。2つのDwordのみが書き込まれる場合、ブリッジは上位32ビットがゼロであると仮定する。ブリッジ106はデバイス・セレクト(DEVSEL#)信号をアサートすることによって、メールボックス・レジスタ216に対するライト処理を横取りし要求する。何故なら、アドレスは二次バス112上のいずれのデバイスにも割り当てられていないからである。ブリッジ106がトランザクションを受け取るとき、バッファが使用可能であればそれを予約し、必要であれば、他の中間ブリッジに送出するために、メールボックス216にアドレス、長さおよびその他の情報を格納する。予約が行われた場合、上述のように、予約情報ブロックの初期設定も行う。したがって、この代替案では、特別アドレスに対するライトが、開始アドレスおよび転送すべきバイト数を示す。この情報は、ブリッジ106が要求元マスタ108のためにデータを先取りするために利用する。

【0038】図5および図6に示すように、ブリッジ階層がある場合、バッファ予約を支持する各ブリッジは、予約情報を上流に渡すことができる。各ブリッジのコンフィギュレーションでは、階層内の最上位のブリッジは、予約を送出することをディスエーブルにされる。図5は、バッファ予約を支持する2つのブリッジ106aと106bの間に挟まれたブリッジ300が非予約ブリッジであるブリッジ階層を示す。この例では、グローバル・メモリ・アドレスに対するライト・トランザクションを用いてブリッジ間で予約情報を受け渡すので、コンフィギュレーション・ビットおよびブリッジ106bがセットされ、予約が非予約ブリッジ300に送出され、予約ブリッジ106bに渡されるようにする。非予約ブリッジ300は、通常のメモリ・トランザクションと予約ライト・トランザクションとの間の相違を認識しない。ブリッジ106aがその一次側でメモリ104に結合されている場合、コンフィギュレーション・ビットおよびブリッジ106aが、メモリへの予約情報の送出をディスエーブルするようにセットされる。このようにして、バッファ予約に対応可能な各ブリッジは、予約情報を受け取る。

【0039】図6は、最上部のブリッジ300はバッファ予約に対応しないが、下位2つのブリッジ106aおよび106bはそれに対応する、ブリッジ階層を示す。この例では、ブリッジ300は最上位ブリッジであり、バッファ予約を支持しないので、ブリッジ106aのコンフィギュレーション・ビットはクリアされ、予約情報

の送出がディスエーブルされる。予約はブリッジ106aと106bの間で受け渡しされる。このように、非予約ブリッジが最上位に位置する場合、非予約ブリッジの下にある予約ブリッジは、予約の送出をディスエーブルされる。最初の代替案に対しての、この方法の利点の1つは、スクリプト・ドリブン・バス・マスタ(script-driven bus master、スクリプト駆動のバス・マスタ)は、殆ど変更せずに、メモリ・ライト・トランザクションを容易に実施できることである。

【0040】双方の代替案は、より長くより効率的なトランザクションにおいて、ブリッジ106がデータが先取り(プリフェッチ)できるようにするものである。バッファ予約を利用することにより、要求されたデータのみが先取りされるので、使用されないデータが先取りされることが回避される。要求されたデータのみが先取りされ、要求元が転送を完了するまで、バス・マスタはこのデータを使用しないと考えることができるので、メモリ104とバッファ206および210に格納されているデータの間のコヒーレンシに関する問題も回避される。このように、好適実施例にしたがって先取りされたデータは、転送が完了するまで保持され且つ有効と見なすことができる。

【0041】遅延トランザクション(delayed transaction)として知られている技法によって、リード動作をブリッジ106によって実行することができる。遅延トランザクションに関する更なる詳細は、上記で援用したPCI仕様及びこの明細書の従来技術の項に記載されている。しかしながら、端的に言えば、遅延トランザクションとは、所定の指定時間内にデータ・トランザクションを完了することができないPCIブリッジおよびPCI I/Oによって用いられる、終了およびリトライ(retry、再試行)の方法である。第1の要求元バス・マスタが低速のI/Oコントローラまたはブリッジが応答するのを待つ代わりに、このI/Oコントローラまたはブリッジは第1の要求元バス・マスタをバスからオフに強制し、待ち状態の第1のバス・マスタが保持することによって通常は浪費されるバス帯域を、他のバス・マスタに使用させる。通常、遅延トランザクションは、バス・マスタによる要求、ターゲットによる要求の完了、およびバス・マスタによるトランザクションの完了、という3段階で完了するように進展する。

【0042】PCIバスの用語では、遅延リード要求とは、発生元のバス上で完了する前に宛先バス上で完了しなければならないトランザクションのことであり、I/Oリード、コンフィギュレーション・リード、メモリ・リード、メモリ・リード・ライン、またはメモリ・リード多重コマンド(memory read multiple command)がそれに該当する可能性がある。図1を参照して、バス・マスタ108がメモリ104からのデータを要求する場合、二次バス112が発生元バスとなり、一次バス110が

宛先バスとなる。一旦宛先バス上で要求が試行されると、これが宛先バス上で完了するまで、連続して繰り返されなければならない。そのときまで、遅延リード要求は単なる1要求に過ぎず、バス・マスタ108はこの要求を後に繰り返さなければならないので、バスのデッドロックを防止するためまたは処理を促進するために、いつでも破棄され得る。遅延リード完了とは、宛先バス上で完了し、完結するために要求元バスに移動しつつあるリード・トランザクションのことである。遅延リード完了は、バス・マスタ108によって要求されたデータ、およびターゲット(メモリ104)のステータスを含む。一旦遅延リード要求が宛先バス上で実行されたなら、これは遅延リード完了となる。

【0043】バッファ予約機構を用いない単純な遅延トランザクションの例は、バッファ領域204上におけるその効果を示すのに役立つであろう。尚、遅延トランザクションは、ここに開示するバッファ予約機構を用いても、または用いなくても使用可能であることに留意されたい。また、バッファ予約機構は、遅延トランザクションを用いても、または用いなくても使用可能であるが、最大の処理性能を得るために、これらを一緒に用いるものと仮定する。

【0044】図1を一例として用いる。第1フェーズの間、バス・マスタ108は、二次バス112上にメモリ・リード・トランザクションを発生する。ブリッジ106は、このトランザクションをデコードし、主メモリが一次バス上にあることを認識し、DEVSEL#信号をアサートすることによってトランザクションを要求し、アクセスを完了するために必要なトランザクション情報をラッチし、従来のPCIリトライ技法を用いて要求を終了する。ラッチされた要求情報のことを、遅延要求、あるいは本例では更に具体的に、遅延リード要求と呼ぶ。PCIの規則によれば、バス・マスタ108はトランザクションを繰り返さなければならない。その間、ラッチされたトランザクション情報を用いて、ブリッジ106は、一次バス110上でメモリ・リード・トランザクションを発生する。ターゲット(主メモリ104)の準備ができていない場合、ブリッジ106にリトライを発行することにより、他の遅延トランザクションを強制することも可能である。しかしながら、メモリ104がリトライを発行しないと仮定すれば、メモリ104は要求を処理し、ブリッジ106は完了ステータスおよび要求されたデータを、そのバッファ204の内の1つに受け取る。遅延要求がライト動作である場合、ブリッジ106は主メモリ104から完了ステータスを取得する。宛先バス(一次バス110)上で遅延要求(遅延された要求)を完了した結果、遅延完了(delayed completion)が生成される。遅延完了は、遅延要求および完了ステータスのラッチされた情報ならびにデータから成る。ブリッジ106は、バス・マスタ108が最初の要求を繰り返すま

で、データおよび完了ステータスをバッファ領域204に格納する。

【0045】第3フェーズの間、バス・マスタは首尾良くバスの再調停を行い、元のメモリ・リード要求を再発行する。ブリッジ106はこの要求をデコードし、処理を要求し、データをバス・マスタ108に供給する。この時点で、遅延完了は引込められ、トランザクションは完了したことになる。バス・マスタ108に返されるステータスは、ブリッジ106が、ターゲット（メモリ104）から、リード・トランザクションを実行したときに取得したものと正確に同一であり、即ち、マスタ・アボート、ターゲット・アボート、パリティ・エラー、ノーマル、またはディスコネクトである。

【0046】遅延トランザクション技法を、ここで開示するバッファ予約機構と組み合わせて用いることにより、バス・マスタとバス・スレーブの間のリード動作は、以前よりも効率的に実行可能となる。

【0047】バッファ領域204の管理は、効率的なデータ転送を得るために望ましいものである。バス領域204の基本的な構成要素は、128バイト・バッファであり、図9のAでは、バッファ500と総称的に呼ばれている。リード先取りデータおよびポストッド・ライト・データ(posted write data)は双方ともブリッジ106（の異なるバッファ）に格納することができるが、リード先取りデータは、ポストッド・メモリ・ライトが完了するまで、有効となることはできない。データがバッファ500に書き込まれると、BUFF_IN_PNTRが増分され、次のアドレスを示す。データは、BUFF_OUT_PNTRによって示されるアドレスに始まって、バッファ500から消去される。BUFF_VALID_PNTRは、バッファ500のどの領域が有効かを示す。BUFF_OUT_PNTRとBUFF_VALID_PNTRの間のアドレスが有効アドレスとなる。遅延リード要求が開始されると、BUFF_IN_PNTRおよびBUFF_OUT_PNTRは、PCIトランザクション・アドレスのビット4〜2に初期設定される。データの各Dwordが到達すると、BUFF_IN_PNTRによって示されるアドレスに格納され、ポインタが増分される。宛先バス上のターゲットに宛てた、未だ実行されていない以前のポストッド・メモリ・ライトがない場合、BUFF_IN_PNTRを増分するときにBUFF_VALIDポインタを増分することができる。前のポストッド・メモリ・ライトがある場合、このライトを完了した後にBUFF_VALID_PNTRを更新する。バス・マスタ108がバッファからデータを取り込む際、各Dwordが除去された後にBUFF_OUT_PNTRが増分される。

【0048】図9のBは、予約資源バッファ206がどのように構成されているかを示す。予約バッファ206は、RSRV_BASE_ADDRによって示される開始アドレス、およびRSRV_LIMIT_ADDRによって示される終了アドレスを有する。予約バッファ206は、3つの128バイト・データ・バッファ（図9のAに示した形式の）で構成さ

れている。リード・トランザクションが進展するに連れて、これら3つのバッファが再使用され、RSRV_LIMIT_ADDRからRSRV_BASE_ADDRを減じた長さの1つの長いバッファをシミュレートする。個々のバッファは、各々、BUFF_STARTによって示される開始アドレス、およびBUFF_LIMITによって示される終了アドレスを有する。

【0049】バッファの状態は、2つの変数、BUFF_STATEおよびPREFETCH_FLAGによって示される。PREFETCH_FLAGは、先取り動作と実際のリード要求との間で区別するために使用される。バッファに格納されているデータの状態を図10に示す。バッファの状態は、システム・リセット、バッファ・タイムアウト、またはバッファ・フラッシュ(buffer flush)の後にはEMPTY（空）と表される。バッファがEMPTYと示されている場合、それを割り当てに使用可能である。バッファがリード要求に使用されている場合、またはリード先取りを実行している場合、このバッファはREQUEST（要求）と表される。一旦あるデータがバッファに到達し始めたなら、状態はREQUESTからPART_COMP（部分完了）に遷移する。一旦全てのデータが受け取られたなら、バッファの状態はPART_COMPからCOMPLETE（完了）に遷移する。バッファの状態がREQUEST、PART_COMPまたはCOMPLETEであり、バッファがフラッシュ（消去、flush）されている場合、状態はEMPTYに戻る。バッファがCOMPLETEと示され、バス・マスタがバッファからのデータ全てを読み取る場合、状態はCOMPLETEからEMPTYに変化する。

【0050】予約アルゴリズムでは、バッファがフラッシュされるのは、データの最後のDwordがバッファから除去されるときである。しかしながら、先取りされたデータのフラッシュは、デフォルト・アルゴリズムまたは保存アルゴリズム(conservative algorithm)では適宜行われ、バス・マスタには決して古いデータを与えないことを保証する。このために、バッファの状態を追跡するのは重要である。

【0051】先取りデータとは、マスタからの実際の遅延リード要求で満たされていたバッファに残されているデータのことを意味するが、マスタが再接続したときには、マスタがディスコネクトする前にデータの一部のみが取り出されている。マスタの先取りデータは、以下の状態の下でフラッシュされる。（1）マスタが、当該マスタの最後のリードとは連続しないアドレスからリードをしようとする場合。（2）CPU102がこのマスタに影響を与え得るライト・トランザクション、例えば、二次バス・ターゲットへのI/Oライト、汎用二次メモリ範囲レジスタへのメモリ・ライト、またはこのマスタのためのスロット特定メモリ範囲レジスタへのメモリ・ライトのようなライト・トランザクションを行う場合。（3）二次バス・マスタがこのマスタに影響を与え得るライト・トランザクション、例えば、二次バス・ターゲットへのI/Oライト、汎用二次メモリ範囲レジスタへ

のメモリ・ライト、このマスタのためのスロット特定メモリ範囲レジスタへのメモリ・ライト、またはいずれかのコンフィギュレーション・レジスタへのCPU102のライトのようなライト・トランザクションを行う場合。

【0052】PREFETCH_FLAGに対する状態について、図11を参照しながら説明する。BUFF_STATEがREQUEST、P

A=(先取り) OR

(マスタ・リード後 AND BUFF_STATE=(PART_COMP OR COMPLETE));

B=(フラッシュ AND BUFF_STATE=(REQUEST, PART_COMP OR COMPLETE));

C=(マスタ・リード・トランザクション AND 遅延リード要求) OR

(先取り完了 AND BUFF_STATE=EMPTY) OR

(フラッシュ AND BUFF_STATE=COMPLETE) OR

(実際の遅延リード要求);

D=(マスタ・リード・トランザクション AND 遅延リード要求) OR

(先取り完了 AND BUFF_STATE=EMPTY);

【0054】バッファ204は、それがEMPTYの場合、またはそれが先取りデータを収容しているが当該バッファに対して進行中の先取りが既にある場合にのみ、新しいREQUESTに対して使用可能となる。バス・マスタ108がリードを行おうとして、ブリッジ106が遅延リード要求を開始した場合、BUFF_STATEはREQUESTに変更され、PREFETCH_FLAGはクリアされる。ブリッジ106の先取りによってリード動作が開始された場合、BUFF_STATEはREQUESTに変更され、PREFETCH_FLAGがセットされる。

【0055】最初のデータが到達したときに、BUFF_STATEはPART_COMPに変更される。マスタがリード・トランザクションを繰り返す、状態がPART_COMPまたはCOMPLETEである場合、リード・データをバス・マスタ108に供給する。バス・マスタ108が、ブリッジ106からディスク接続する前に、全てのデータを取り込まなかった場合、PREFETCH_FLAGがセットされ、バッファ204内のデータの残りはこの時点では先取りされたデータと見なされることを示すが、状態は不変のままである。バス・マスタ108が最後のデータを取り込んだ場合、即ち、状態がCOMPLETEで(BUFF_OUT_PNTR >= BUFF_IN_PNTR); リード後)の場合、状態はEMPTYに変更される。

【0056】バッファ204に対してフラッシュ・イベントが発生し、そのPREFETCH_FLAGがセットされ、状態がREQUEST、PART_COMPまたはEMPTYと示されている場合、先取りが現在進行中であるので、状態はEMPTYに変更されるが、未処理の先取りが完了するまで、PREFETCH_FLAGはセットされたままである。バッファ204は、PREFETCH_FLAGがクリアされるまで、他の先取りのために再使用することはできない。バッファ204に対してフラッシュ・イベントが発生し、そのPREFETCH_FLAGがセットされ、状態がCOMPLETEと示された場合、状態はEMPTYに変更され、PREFETCH_FLAGはクリアされる。

【0057】PREFETCH_FLAGがセットされている間に、

ART_COMPまたはCOMPLETEと示された場合、PREFETCH_FLAGは、リードが先取りか又は実際のREQUESTかを示す。PREFETCH_FLAGがセットされている場合、先取り動作が進行中であり、一方PREFETCH_FLAGがクリアされている場合、実際のREQUESTが進行中である。セット状態からクリア状態への遷移は、以下の式で示される。

【0053】

実際の遅延要求がバッファ204を必要とする場合、新しい要求がバッファに与えられる。PREFETCH_FLAGはクリアされ(入来する先取りデータを破棄させる)、状態はREQUESTに変更される。

【0058】先取りデータがバッファ204に到達し、PREFETCH_FLAGがクリアされると、データは破棄される(バッファは実際の要求のために再使用されている)。先取りデータがバッファ204のために到達し、状態がEMPTYの場合、データは破棄される(バッファがフラッシュされる)。最後の先取りデータが到達したとき、PREFETCH_FLAGがクリアされ、バッファ204は使用可能となり再使用される。

【0059】ブリッジ106には、3つの先取りおよびフラッシュ・アルゴリズム、即ち、デフォルト・アルゴリズム、保存アルゴリズム、および予約アルゴリズムが実施されている。デフォルト・アルゴリズムは、殆どのバス・マスタに適用可能でなければならない。保存アルゴリズムは、先取りを少なく、フラッシュを多くするが、デフォルト・アルゴリズムが特定のバス・マスタに対して作用しない場合にのみ実施される。予約アルゴリズムは、予約機構を開始可能な、一層高度なバス・マスタによって用いられる。

【0060】次に図12ないし図19を参照しながら、好適実施例による先取り(プリフェッチ)およびフラッシュ・アルゴリズムについて説明する。ブリッジ106は、データのリードおよびライトに関する種々の処理を実行するが、簡略化のために、これらのステップを省略してあることは理解されたい。更に、ライト・ポストイングに関連するステップも示されていないが、ブリッジは、ここに記載するリード先取りおよびフラッシュ・プロセスと共に、これらの処理も実行することは理解されよう。

【0061】先取りおよびフラッシュ・アルゴリズムは、図12において、二次バス・インターフェース20

2上のバス・トランザクションの受け取りから開始する。ステップ700(図12)において、コンフィギュレーション・コマンドが受け取られたか否かについて判定が行われる。受け取られた場合、制御はステップ710(図13)に移行し、このコマンドを完了する。その他の場合、制御はステップ702に進む。ステップ702において、ライト動作がブリッジ106に発行されたか否かについて判定が行われる。発行された場合、制御はステップ712(図14)に移行し、ライト動作が、既に割り当てられているデフォルト・バッファの1つに影響を与えるか否かについて判定を行う。バス・マスタ108が、デフォルト・バッファ210の内の割り当てられた1つであり、ライト動作が、コンフィギュレーション・ライトまたはコマンドのように、バス・マスタ108に対して実行される場合、新しいプロセスが開始された可能性があり、この場合、デフォルト・バッファ210内にある先取りされたデータはいずれもが古い又は一致しない恐れがある。かかる場合、バッファはライト・コマンドによる影響を受けたと言う。このような場合、制御はステップ714に移行し、ライト動作によって影響を受けたスロット又はマスタに対応するデフォルト・バッファ210をフラッシュする。次に、制御はステップ714からステップ716に移行する。ステップ712において、ライト・コマンドがいずれのデフォルト・バッファ210にも影響を与えないと判定された場合、制御はステップ712からステップ716に移行する。ステップ716において、あらゆるライト・データのライト・ポスティング・バッファへの受け取り及びBUFF_VALID_PNTRの更新(ステップ718)を含む、残りのライト動作処理が行われる。BUFF_VALID_PNTRは、デフォルト・アルゴリズムおよび保存アルゴリズムについてのみ更新される。これで、このルーチンは終了し、それをコールしたいいずれかのプロセスに戻る。

【0062】ステップ702に戻り、動作がライト動作でないと判定された場合、制御はステップ704に移行し、動作が予約コマンドであるか否かについて判定が行われる。予約コマンドである場合、制御はステップ720(図15)に移行し、予約コマンドがブリッジ106によって受け取られる。次に、制御はステップ722に移行し、予約資源220の1つがこの要求のために使用可能であるか否かについて判定が行われる。使用可能ではない場合、制御はステップ742に移行し、要求元バス・マスタにリトライを発行し、このルーチンは終了する。予約資源220が使用可能である場合、制御はステップ722からステップ724に移行する代わりに、予約コマンド内に含まれているデータ(即ち、アドレスおよびカウンタ)を用いて、予約資源レジスタを初期設定する。次に、制御はステップ728に移行し、3つの128バイト・バッファをこの要求に割り当て、ステップ730に移行し、これらのバッファに対するステータス

をEMPTYに変更する。次に、制御はステップ730からステップ732に移行し、ブリッジ106はデータの3ないし4ラインを、バッファ206の第1の128バイト・バッファに先取りする。次に、制御はステップ734に移行し、最後のバッファ・ポインタ(LAST_BUFF_PNTR)を初期設定し、どのバッファが次に書き込まれるのかを示す。

【0063】次に、制御はステップ736に移行し、送出コンフィギュレーション・ビット(forwarding configuration bit)がセットされているか否かについて判定が行われる。セットされている場合、次に制御はステップ738に移行し、予約コマンドをいずれかの上流ブリッジに向けて送出する。次に、制御はステップ738からステップ740に移行する。ステップ736において、コンフィギュレーション・ビットがセットされていない場合、制御は直接ステップ740に移行する。ステップ740において、予約タイマをリセットし、開始する。25ms、50msまたは100msのような所定の時間量の後、予約資源にアクセスされなかった場合、予約タイマはこの予約を取り消す。この無動作期間の後、データは古くなったものと見なす。これで、このルーチンは終了する。

【0064】ステップ704に戻り、動作が予約コマンドではないと判定された場合、次に制御はステップ706に移行し、動作はバス・マスタからのリード・リトライか否かについて判定が行われる。バス・マスタからのリード・リトライでない場合、この動作は最初のリード動作であると見なされ、制御はステップ744(図17)に移行する。

【0065】ステップ744において、要求元バス・マスタが、それにデフォルト・バッファ210を既に割り当ててあるか、または新しいバス・マスタであるかについて判定が行われる。これが新しいマスタでないと判定された場合、制御はステップ746に移行し、要求されたアドレスが、BUFF_OUT_PNTRによって指し示されているアドレスに等しいか否かについて判定が行われる。等しい場合、制御は次にステップ750に移行する。等しくない場合、制御は次にステップ748に移行し、マスタは既に、新しい、連続しないアドレスからリードを開始しているので、割り当てられたデフォルト・バッファ210をフラッシュする。次に、制御はステップ750に移行し、デフォルト・バッファ210をこのマスタに再度割り当て、上述のように、バッファ・レジスタを初期設定する。次いで、制御はステップ768に移行する。

【0066】ステップ744において、これが新しいマスタであると判定された場合、制御はステップ756に移行し、デフォルト・バッファ210のいずれかが現在EMPTY(空)と示されているか否かについて判定が行われる。EMPTYと示されているデフォルト・バッファがあ

る場合、制御はステップ758に移行し、EMPTYのデフォルト・バッファ210の1つを新しいマスタに割り当てる。ステップ758から、制御はステップ768に移行する。EMPTYと示されているデフォルト・バッファがない場合、制御はステップ760に移行し、いずれかのデフォルト・バッファ210に、現在REQUESTまたはPART_COMPと示されているものがあるか否かについて判定が行われる。かかるバッファがある場合、制御はステップ726に移行し、PREFETCH_FLAGがセットされているか否かについて判定が行われ、セットされている場合、制御はステップ764に移行し、デフォルト・バッファ210が新しいマスタに割り当てられる。次に、制御はステップ764からステップ768に移行する。ステップ760および762において、判定が否定的である場合、制御はステップ766に移行し、バス・マスタに、後の時点で動作をリトライするように要求する。

【0067】ステップ768において、制御は、新しく割り当てられたまたは再割り当てされたデフォルト・バッファ210に対する先取りをディスエーブルする。次に、制御はステップ770(図16)に移行し、リード要求情報を格納する。次に、制御はステップ772に移行し、ブリッジ106が宛先バス上のターゲットから要求されたデータをcollectingしている間に、バス・マスタに動作をリトライするように要求する。次に、制御はステップ774に移行し、リード要求がメモリ・リード多重コマンド(memory read multiple command)であったか否かについて判定を行う。そうであった場合、制御はステップ776に移行し、ブリッジ106は、主メモリ104からデータの3ないし4ラインを読み取る。次に、制御はステップ776からステップ784に移行する。ステップ774において、リード・コマンドがメモリ・リード多重でなかったと判定された場合、制御はステップ778に移行し、リード・コマンドはメモリ・リード・ライン・コマンドであったか否かについて判定を行う。そうであった場合、制御はステップ780に移行し、ブリッジ106が主メモリ104からデータの1行を読み取り、次いで、制御はステップ784に移行する。ステップ778において、コマンドがメモリ・リード・ライン・コマンドでなかったと判定された場合、制御はステップ782に移行し、ブリッジ106は主メモリ104からデータの1Dwordを読み取り、次いで、制御はステップ784に移行する。

【0068】ステップ784において、ブリッジ106は主メモリ104からデータを受け取り、ステップ786において、ブリッジ106は受け取ったデータを割り当てられたバッファに格納し、処理は終了する。このルーチンは後に要求元バス・マスタがこの動作のリトライを行うときにコールされる。

【0069】ここでステップ706に戻り、動作がバス・マスタからのリード・リトライ動作であると判定され

た場合、制御はステップ708に移行し、更に、このバス・マスタは以前にバッファ予約を要求したことがあるか否かについて判定が行われる。要求したことがある場合、制御はステップ788(図18)に移行し、このリトライ・リード動作に対するアドレスが、予約資源220の制限内であるか否かについて判定が行われる。制限内でない場合、制御はステップ790に移行して予約を取り消し、制御はステップ744(図17)に移行し、デフォルト・アルゴリズムを用いてリード動作を開始する。ステップ788において、リード動作が予約制限内であることが判定された場合、制御はステップ792に移行し、次に、リード・アドレスが3つのバッファの制限の1つの中にあるか否かについて判定が行われる。そうである場合、制御はステップ794に移行し、要求されたデータのDwordをバス・マスタに供給する。ステップ794から、制御はステップ796に移行し、バス・マスタがこの特定の128バイト・バッファからリードを行ったのは、これが最初であるか否かについて判定が行われる。

【0070】これが最初である場合、制御はステップ798に移行し、ブリッジ106が、LAST_BUFF_PNTRによって示される次のバッファにデータを先取りする。バッファ206は、リセント・リーセントリー・ユーズド(least recently used)(LRU)アルゴリズムにしたがって充填される。一旦あるバッファからデータが最初に読み出されたなら、この次のデータ先取り動作には、最も長い間使用されていないバッファが用いられる。したがって、この点において、バス・マスタは1つのバッファからリードを行い、次のバッファは先取りされたデータを保持し、別のバッファが以前に先取りされたデータを保持する。したがって、バス・マスタがバックアップする必要があるがあっても、あるいはその現アドレスからステップ・バックする必要があるがあっても、データは、以前に先取りされたバッファにある可能性が未だに高い。制御はステップ798からステップ800に移行する。これが最初ではない場合、制御はステップ796からステップ800に移行し、バス・マスタによるリード動作がバースト・リードであるか否かについて判定が行われる。バースト・リードである場合、制御はステップ788に戻り、追加データを供給する。バースト・リードでない場合、制御は終了に移行する。

【0071】ステップ792において、リード動作がバッファ制限内にないことが判定された場合、制御はステップ804に移行し、ブリッジ106は、ブリッジが要求されたデータを取り出している間に、マスタに動作をリトライすることを強制する。次に、制御はステップ804から806に移行し、ブリッジが主メモリ104からデータを読み取り、ステップ808に移行してそのデータを受け取り、ステップ810に移行して受け取ったデータをバッファ206に格納し、ブリッジはバス・マ

スタがリード動作をリトライするのを待つ。

【0072】ステップ708に戻り、このリード動作のために以前に要求した予約資源がなかった場合、制御はステップ812のデフォルト・アルゴリズム(図19)に移行する。

【0073】ステップ812において、デフォルト・バッファ210の1つから、要求されたデータをバス・マスタに供給する。次いで、制御はステップ814に移行し、リード動作がバースト動作であるか否かについて判定が行われる。バースト動作でない場合、処理はステップ815に進み、保存(conservative)アルゴリズムが選択されるか否かについて判定が行われる。される場合、制御はステップ817に移行し、バッファをフラッシュする。保存アルゴリズムが選択されない場合、制御は終了に移行する。ステップ814において、バースト動作が進行中である場合、制御はステップ816に移行し、次の要求されたアドレスが、データ・バッファ内に含まれているか否かについて判定が行われる。そうである場合、制御はステップ818に移行し、データをバス・マスタに供給する。次に、制御はステップ814に戻り、要求されれば、追加データを供給する。ステップ816において、次の要求されたアドレスがバッファ内に含まれていないと判定された場合、制御はステップ820に移行し、このバッファに対して先取りがイネーブルされているか否かについて判定が行われる。イネーブルされていない場合、制御はステップ832に移行し、先取りをイネーブルする。次いで、上述のように、制御はステップ770に移行する。

【0074】ステップ820において、先取りがイネーブルされていると判定された場合、制御はステップ822に移行し、リード・トランザクション情報をブリッジ106に格納する。次に、制御はステップ824および826に移行し、ブリッジ106がバス・マスタにリトライを発行し、ブリッジはリード動作を主メモリ104に発行する。この点以降、このバス・マスタに対して先取りがイネーブルされる。ステップ828において、ブリッジ106はリード・データを受け取り、ステップ830において、ブリッジは受け取ったデータをバッファに格納し、バス・マスタがリード動作をリトライするのを待つ。

【0075】以上の本発明の開示および記載は例示であり、その説明に供するものであり、本発明の精神から逸脱することなく、サイズ、形状、材質、構成要素、回路素子、配線接続および接点、ならびに図示した回路や構成、更に動作方法における種々の変更が可能である。

【図面の簡単な説明】

【図1】好適実施例に従ったブリッジ・デバイスを組み込んだコンピュータ・システムのブロック図である。

【図2】好適実施例に従ったブリッジ・デバイスを組み込んだコンピュータ・システムの更に詳細なブロック図である。

【図3】好適実施例に従ったブリッジ・デバイスの機能ブロックを示すブロック図である。

【図4】好適実施例に従ったブリッジ・デバイスの機能ブロックを示すブロック図である。

【図5】好適実施例に従った2つの予約ブリッジ間に結合された非予約ブリッジを有する、バスの階層の一例を示すブロック図である。

【図6】非予約ブリッジが階層の最上部のホスト・ブリッジとして動作する、バスの階層の一例を示すブロック図である。

【図7】予約コマンドを組み込んだ、PCIバス・リード・トランザクションを示すタイミング図である。

【図8】予約コマンドを組み込んだ64ビット・アドレスによる、PCIバス・リード・トランザクションを示すタイミング図である。

【図9】バッファのブロック図であり、Aは好適実施例に従った単一の128バイト・データ・バッファのブロック図であり、Bは好適実施例に従った予約バッファのブロック図である。

【図10】好適実施例に従ったバッファ状態を示す状態図である。

【図11】好適実施例に従ったバッファ先取り状態を示す状態図である。

【図12】好適実施例に従ったプリフェッチ及びフラッシュ・アルゴリズムを示すフローチャートである。

【図13】好適実施例に従った、図12の「A」に続くフローチャートである。

【図14】好適実施例に従った、図12の「B」に続くフローチャートである。

【図15】好適実施例に従った、図12の「C」に続くフローチャートである。

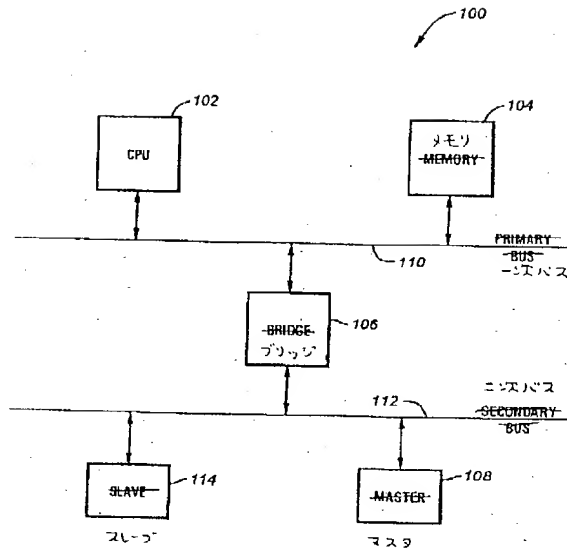
【図16】好適実施例に従った、図17および図19の「G」に続くフローチャートである。

【図17】好適実施例に従った、図12の「D」に続くフローチャートである。

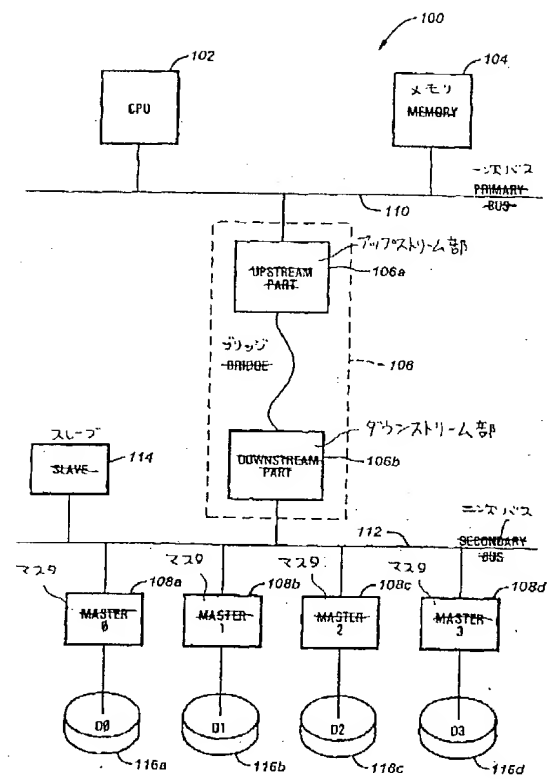
【図18】好適実施例に従った、図12の「E」に続くフローチャートである。

【図19】好適実施例に従った、図12の「F」に続くフローチャートである。

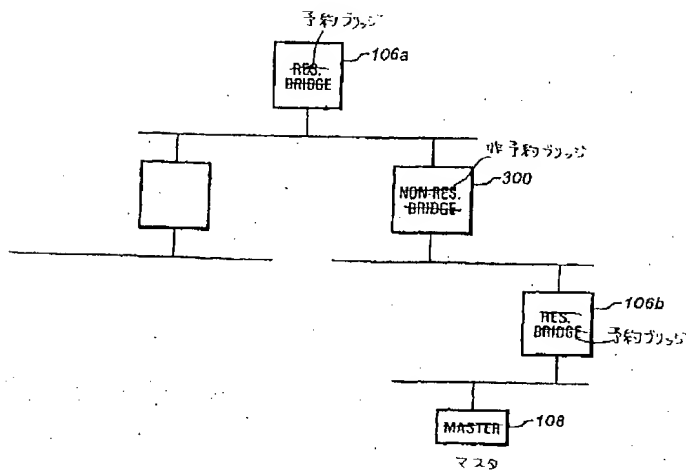
【図1】



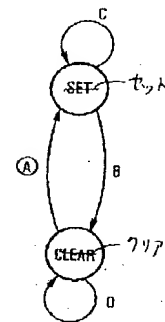
【図2】



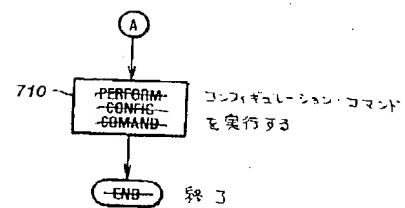
【図5】



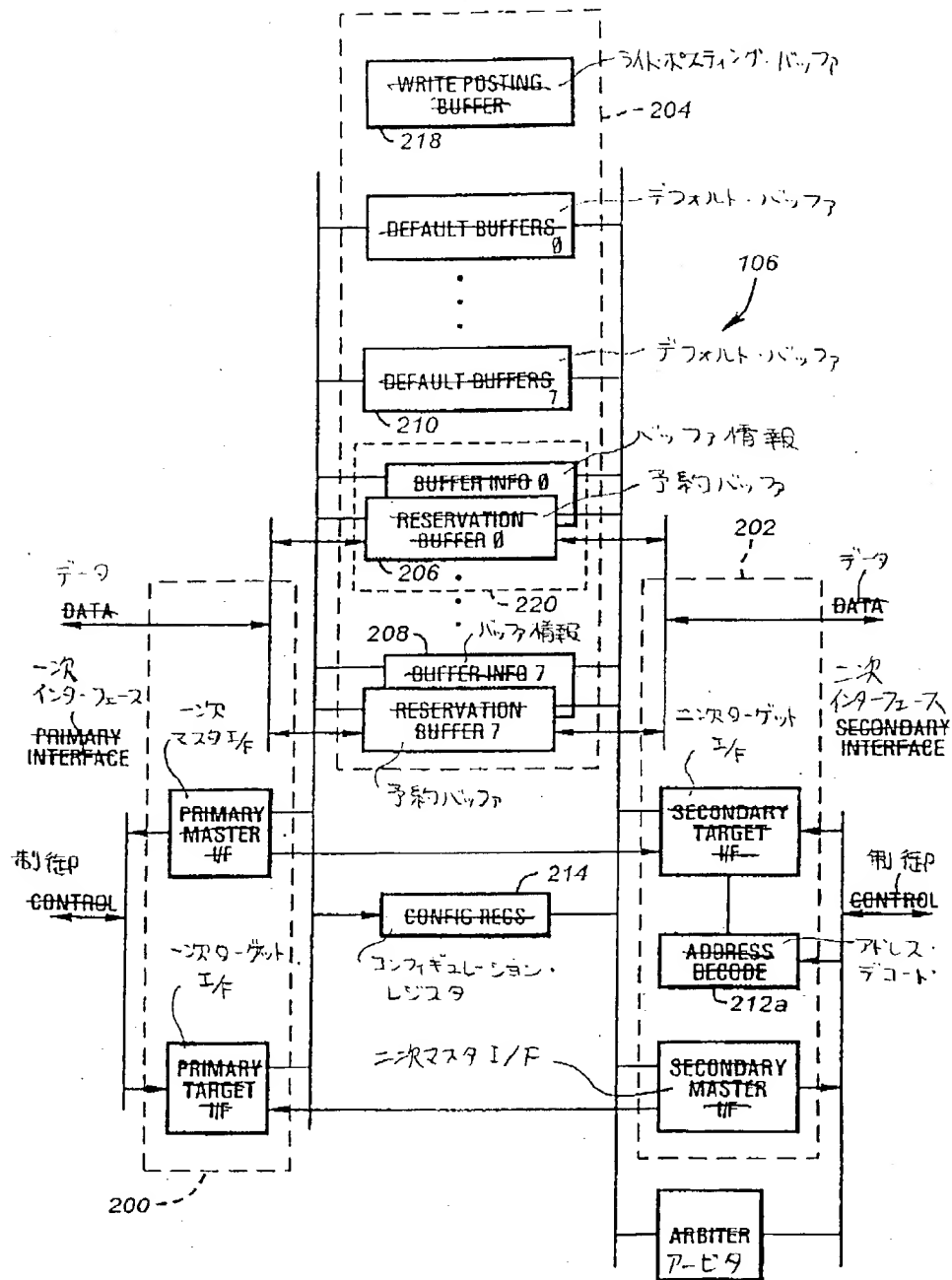
【図11】



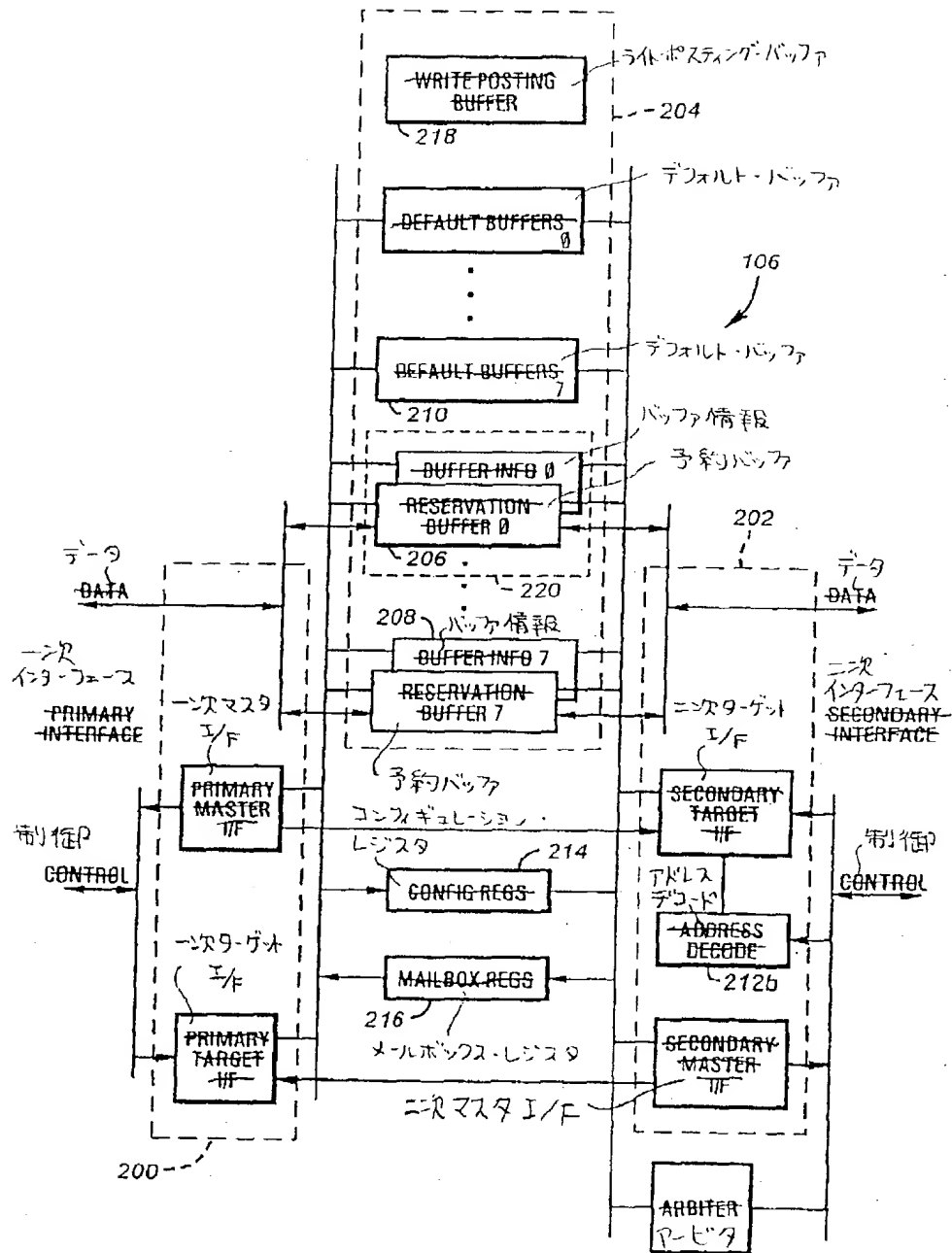
【図13】



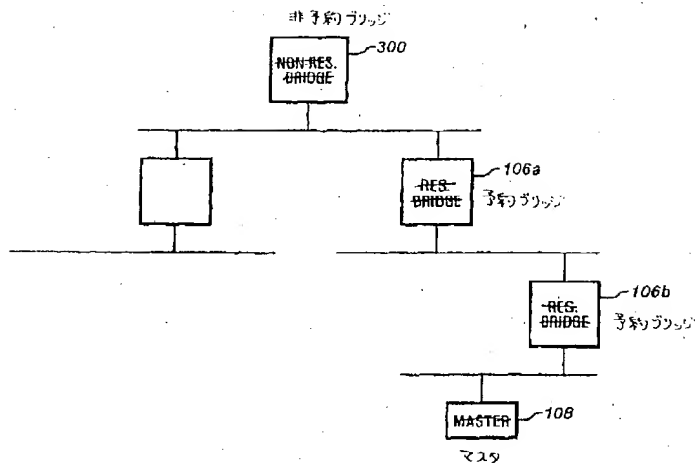
【図3】



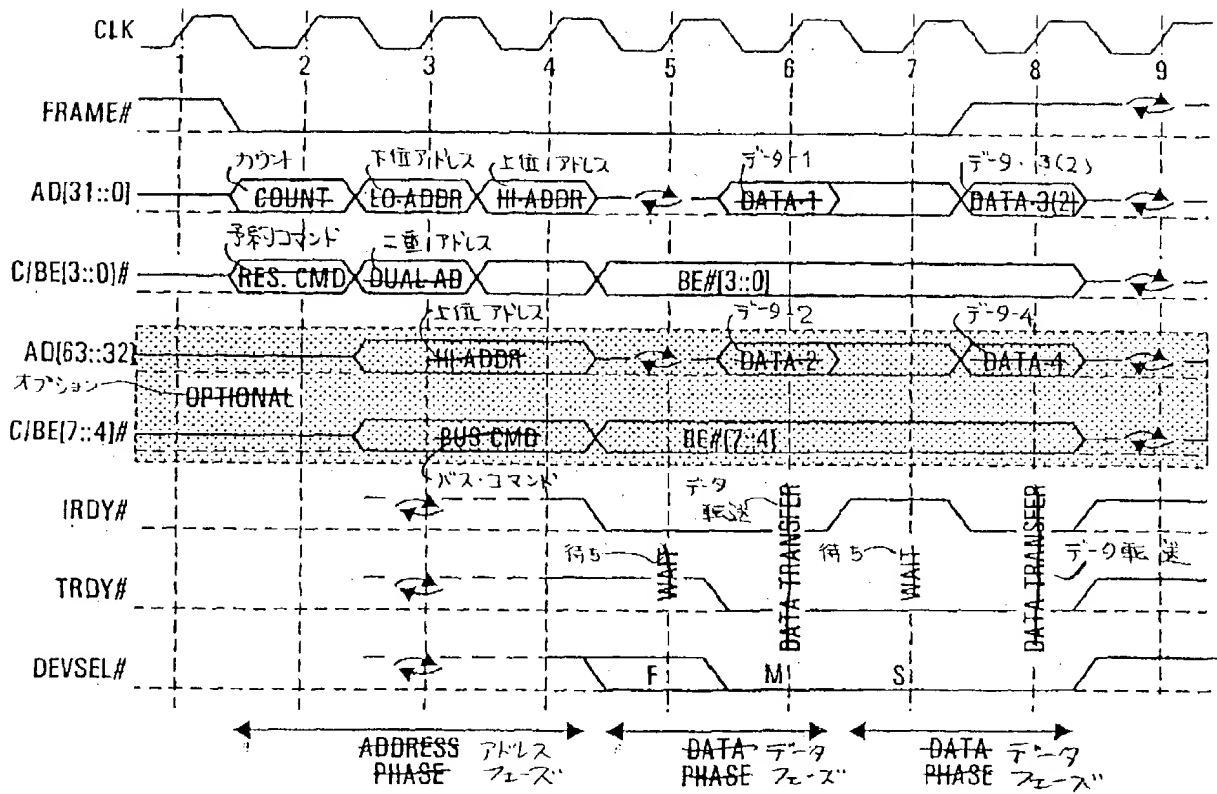
【図4】



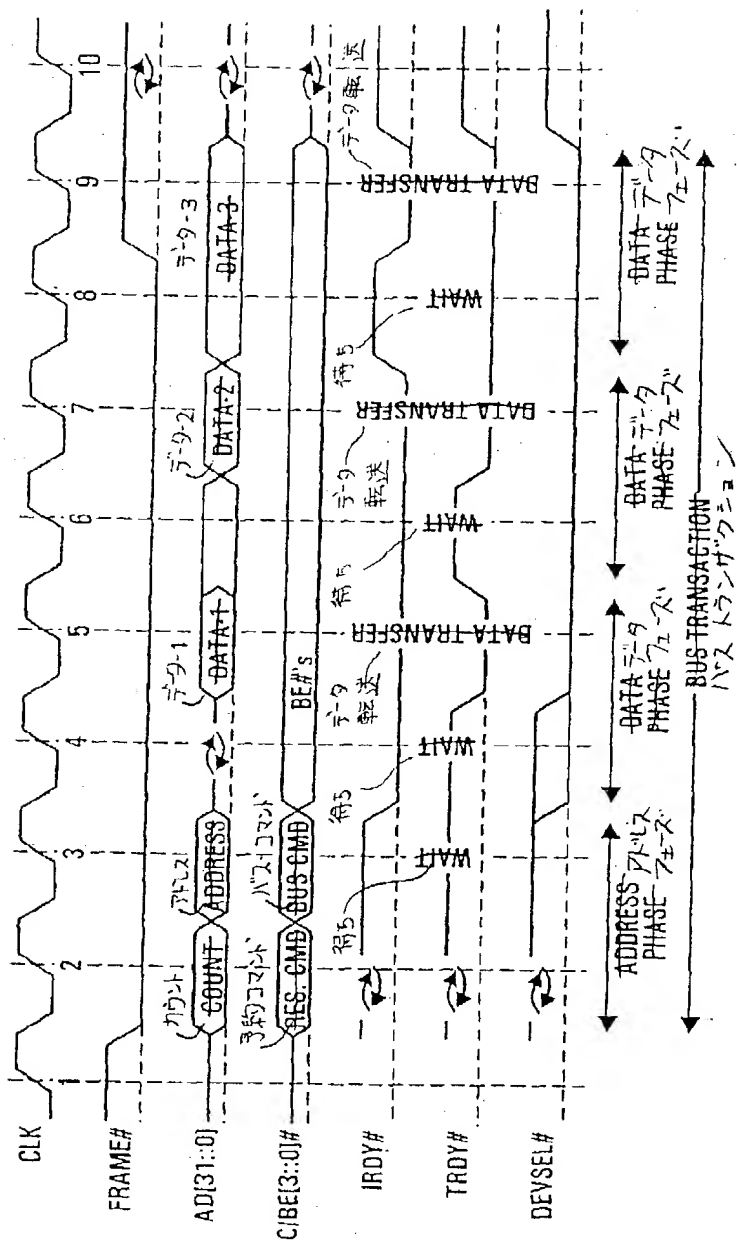
【図6】



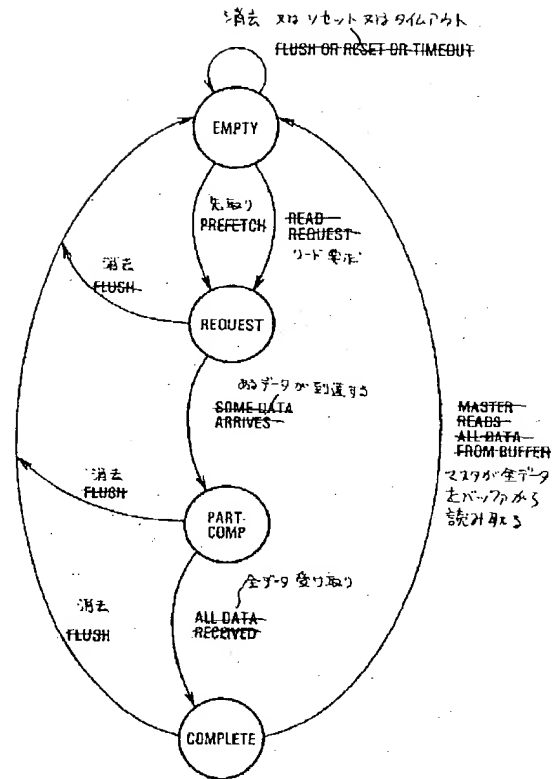
【図7】



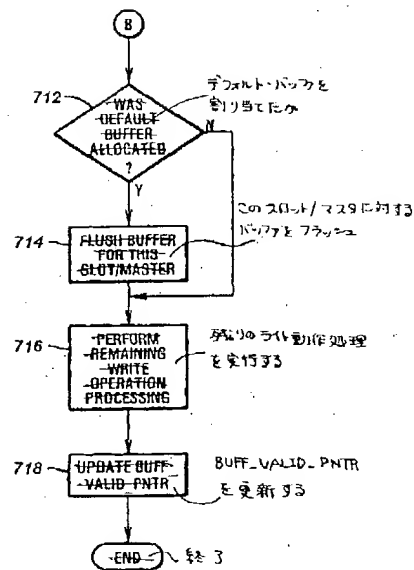
【図8】



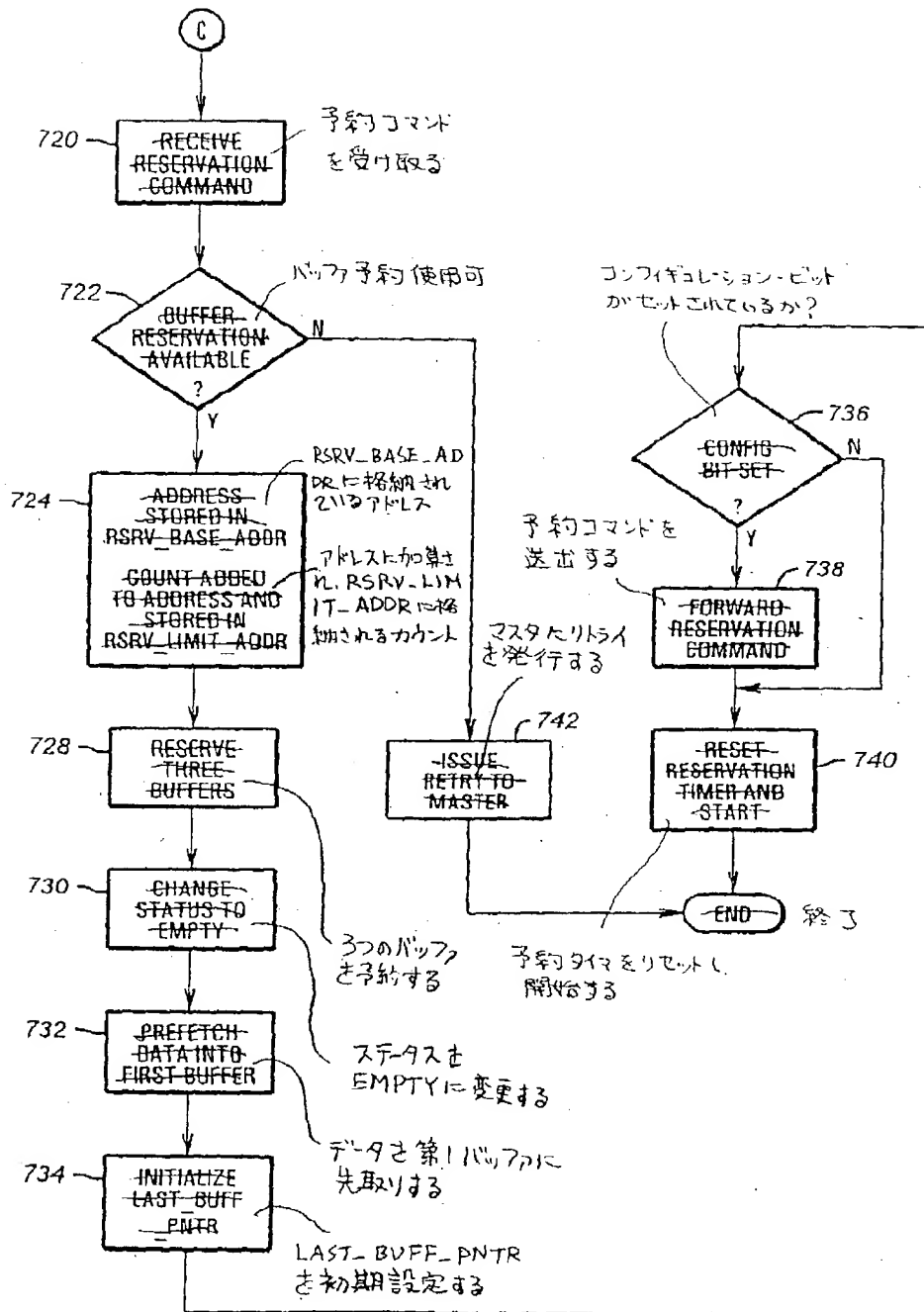
【図 10】



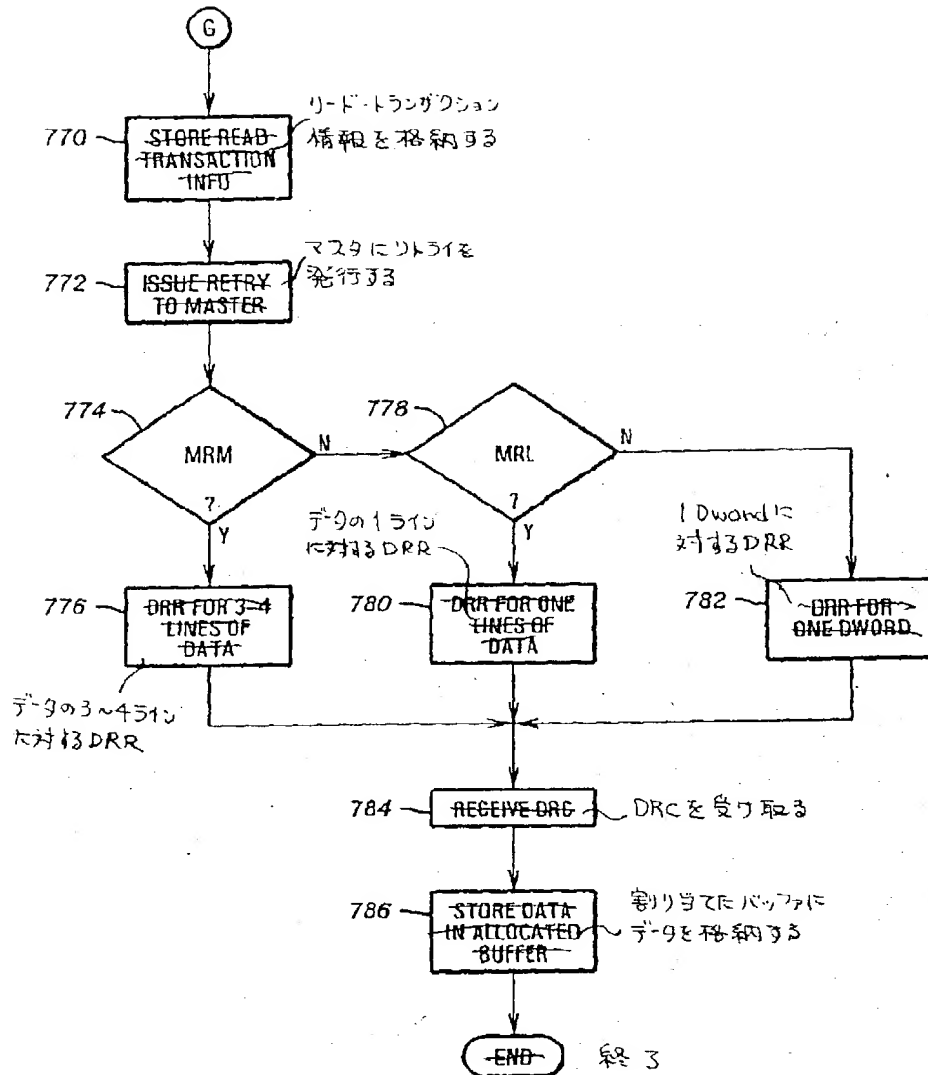
【图 14】



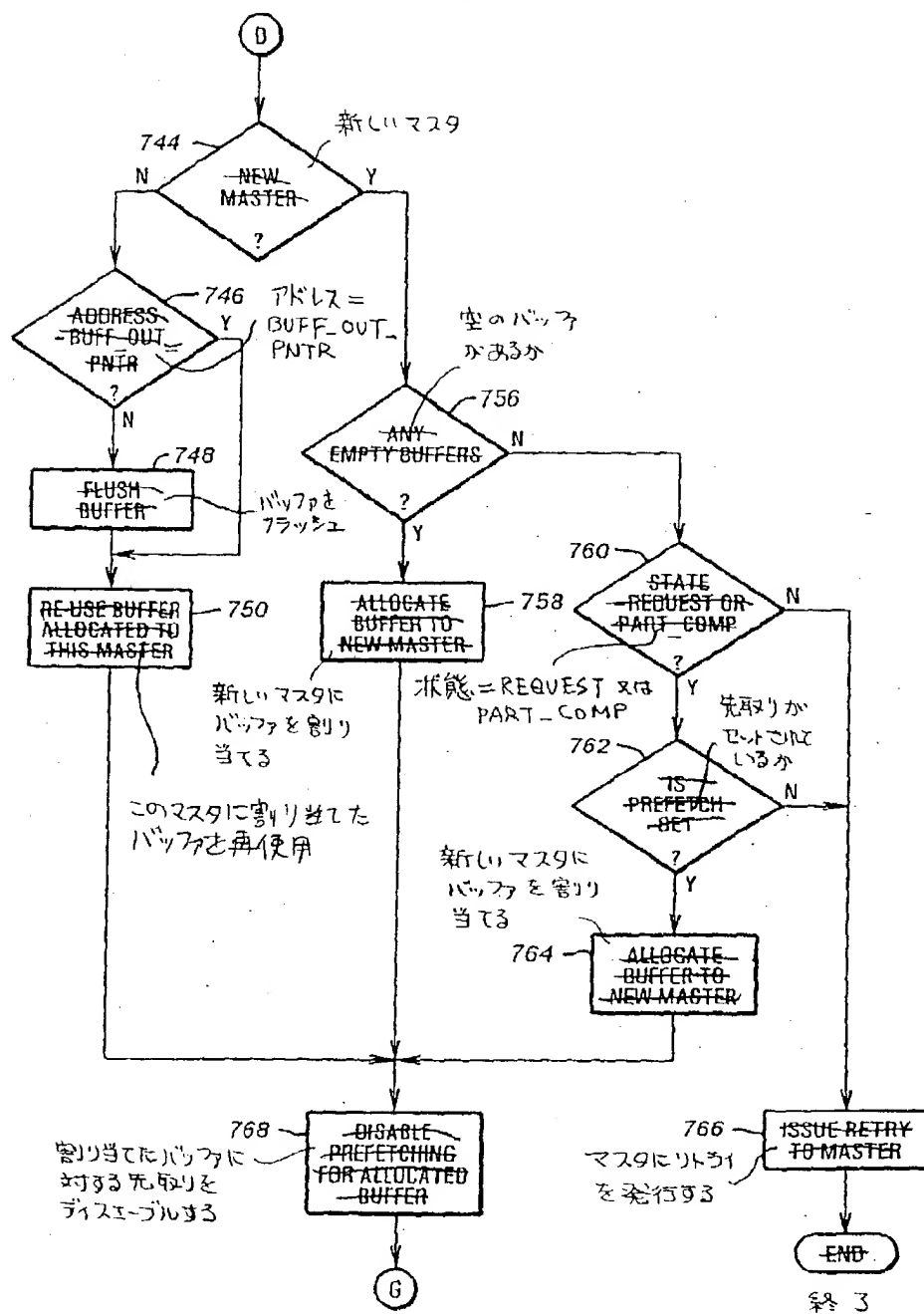
【図15】



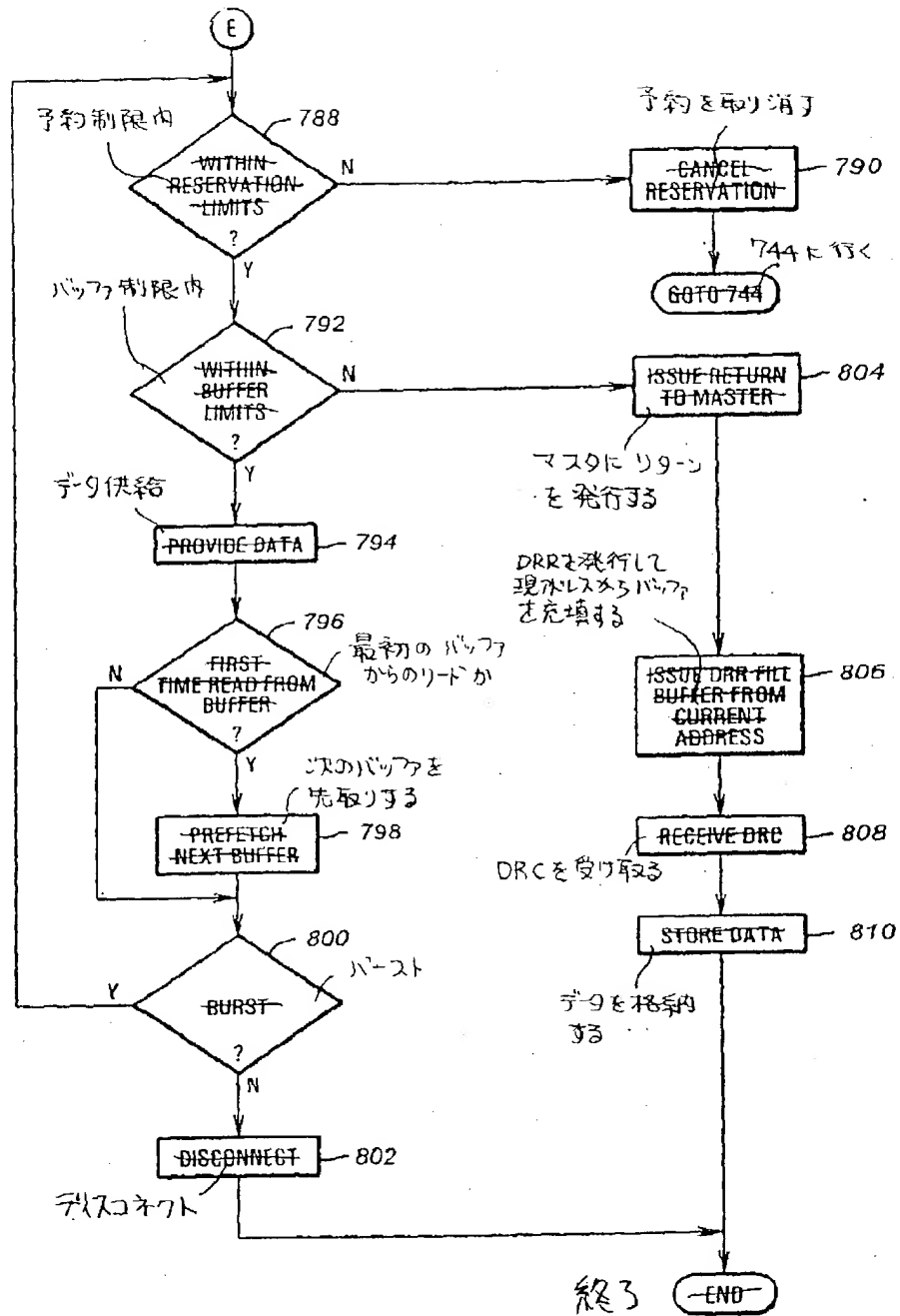
【図16】



【図17】



【図18】



```

graph TD
    F((F)) --> 812[PROVIDE DATA TO MASTER]
    812 --> 814{BURST?}
    814 -- Y --> 816{NEXT ADDR WITHIN BUFFER?}
    814 -- N --> 815{CONSERVATIVE?}
    816 -- Y --> 818[PROVIDE NEXT DATA]
    816 -- N --> 815
    818 --> 814
    815 -- Y --> 817[FLUSH]
    815 -- N --> END1((END))
    817 --> END1
    815 -- 保存 --> 815
    815 -- リード・トランザクション情報を格納する --> 822[STORE READ TRANSACTION INFO]
    822 --> 824[ISSUE RETRY TO MASTER]
    824 --> 826[ISSUE DRB TO TARGET TO PREFETCH BUFFER]
    826 --> 828[RECEIVE DRC FROM TARGET]
    828 --> 830[STORE RECEIVED DATA]
    830 --> END1
    815 -- 先取りがイネーブルされたか --> 820{PREFETCHING ENABLED?}
    820 -- Y --> 822
    820 -- N --> 832[ENABLE PREFETCHING]
    832 --> G((G))
    832 -- 先取りをイネーブルする --> G
  
```

Flowchart illustrating the data transfer process:

- Start (F)**: Master provides data (マスターにデータを供給する).
- 812 PROVIDE DATA TO MASTER**: Data is provided to the master.
- 814 BURST?**: Decision on burst mode.
 - Y**: Proceeds to 816.
 - N**: Proceeds to 815.
- 816 NEXT ADDR WITHIN BUFFER?**: Decision on next address within buffer.
 - Y**: Proceeds to 818.
 - N**: Proceeds to 815.
- 818 PROVIDE NEXT DATA**: Next data is provided.
- 815 CONSERVATIVE?**: Decision on conservative mode.
 - Y**: Proceeds to 817.
 - N**: Proceeds to END.
 - 保存** (Save): Loop back to 815.
 - リード・トランザクション情報を格納する** (Store read transaction information): Proceeds to 822.
 - 先取りがイネーブルされたか** (Is prefetching enabled?): Proceeds to 820.
- 817 FLUSH**: Flush operation.
- 820 PREFETCHING ENABLED?**: Decision on prefetching.
 - Y**: Proceeds to 822.
 - N**: Proceeds to 832.
- 822 STORE READ TRANSACTION INFO**: Store read transaction information.
- 824 ISSUE RETRY TO MASTER**: Issue retry to master.
- 826 ISSUE DRB TO TARGET TO PREFETCH BUFFER**: Issue DRB to target to prefetch buffer.
- 828 RECEIVE DRC FROM TARGET**: Receive DRC from target.
- 830 STORE RECEIVED DATA**: Store received data.
- 832 ENABLE PREFETCHING**: Enable prefetching.
- G**: Proceeds to END.
- 先取りをイネーブルする** (Enable prefetching): Proceeds to G.
- 終了** (End): Process ends.

(71)出願人 591030868
20555 State Highway
249, Houston, Texas
77070, United States o
f America

【 外国語明細書 】

1. Title of Invention

Buffer Reservation Method for a Bus Bridge System

2. Claims

1. A method of communication between a bus bridge and a bus master, the bus bridge having at least one data buffer, the method comprising the steps of:
 - (a) receiving a reservation command from the bus master; and
 - (b) reserving a data buffer for the bus master in response to said reservation command.
2. The method of claim 1, further comprising the step of:
 - (c) receiving a bus command after step (a).
3. The method of claim 2, wherein said bus command includes a starting address and wherein said bus command and said reservation command together define an address range.
4. The method of claim 3, wherein said reservation command includes a count.
5. The method of claim 3, wherein if said bus command is a read operation the method further comprises the steps of:
 - (d) fetching data according to the read operation, the data being passed through the reserved buffer;
 - (e) prefetching data according to the read operation and said address range, the prefetched data being written into the reserved data buffer; and
 - (f) providing the fetched and prefetched data to the bus master.
6. The method of claim 5, further comprising the steps of:
 - (g) repeating steps (e) and (f) until said address range is exhausted; and
 - (h) canceling the data buffer reservation after step (g).
7. The method of claim 1, wherein said reservation command is received after a frame signal is asserted.

8. The method of claim 1, wherein the bus bridge and the bus master are for coupling to a peripheral component interconnect (PCI) bus.
9. A method of communicating between a first and a second bus bridge, the second bus bridge for coupling between the first bus bridge and a bus master, the method comprising the steps of:
- (a) receiving a reservation command at the second bus bridge from the bus master; and
 - (b) forwarding said reservation command to the first bus bridge.
10. The method of claim 9, wherein the second bus bridge includes a programmable forwarding bit to select whether reservation commands are forwarded, the method further comprising the step of:
- (c) determining whether the forwarding bit is set to enable forwarding of reservation commands before step (b), and
- wherein step (b) is conditioned on the forwarding bit being set to enable forwarding.
11. The method of claim 9, the method further comprising the steps of:
- (d) receiving a bus command after step (a); and
 - (e) reserving a data buffer for the bus master in response to said reservation command and said bus command.
12. The method of claim 11, wherein said bus command includes a starting address and wherein said bus command and said reservation command together define an address range.
13. The method of claim 12, wherein said reservation command includes a count.

14. The method of claim 12, wherein if said bus command is a read operation the method further comprises the steps of:

(f) fetching data according to the read operation, the data being passed through the reserved buffer;

(g) prefetching data according to the read operation and said address range, the prefetched data being written into the reserved data buffer; and

(h) providing the fetched and prefetched data to the bus master.

15. The method of claim 14, further comprising the steps of:

(i) repeating steps (g) and (h) until said address range is exhausted; and

(j) canceling the data buffer reservation after step (i).

16. The method of claim 9, wherein said reservation command is received after a frame signal is asserted.

17. The method of claim 9, wherein the bus bridge and the bus master are for coupling to a peripheral component interconnect (PCI) bus.

18. A method of reserving a buffer in a bus bridge, the method comprising the steps of:

(a) providing a frame indication;

(b) providing a reservation command when said frame indication is provided; and

(c) providing a bus command after step (b) and while said frame indication is provided.

19. The method of claim 18, wherein said bus command includes a starting address and wherein said bus command and said reservation command together define an address range.

20. The method of claim 18, wherein said reservation command includes a count.

21. A computer system, comprising:
- a main memory;
 - a processor coupled to said main memory;
 - a mass storage system;
 - a bus master coupled to said mass storage system, said bus master operable to provide a bus operation including a reservation command and a bus command; and
 - a bus bridge coupled between said main memory and said bus master, said bus bridge having at least one assignable data buffer, said bus bridge operable to receive the bus operation and reserve a data buffer for the exclusive use of said bus master in response to said bus operation.
22. The computer system of claim 21, wherein said bus command includes a starting address and wherein said bus command and said reservation command together define an address range.
23. The computer system of claim 22, wherein said reservation command includes a count.
24. The computer system of claim 22, wherein if said bus command is a read operation
- said bus bridge fetches data according to the read operation and prefetches data according to said address range, the data being provided to the bus master when requested.
25. The computer system of claim 24, wherein when said bus master reads at a last address of the address range, said bus bridge cancels the data buffer reservation.
26. The computer system of claim 21, wherein said reservation command is received after a frame signal is asserted.
27. The computer system of claim 21, wherein the bus bridge and the bus master are coupled to a peripheral component interconnect (PCI) bus.

28. The computer system of claim 21, wherein said bus operation includes a memory read operation for writing to said mass storage system.

29. A computer system, comprising:
a main memory;
a processor coupled to said main memory;
a mass storage system;
a bus master coupled to said mass storage system, said bus master operable to provide a bus operation including a reservation command and a bus command;
a first bus bridge coupled to said main memory; and
a second bus bridge coupled between said first bus bridge and said bus master, said second bus bridge having at least one assignable data buffer and a forwarding bit, said second bus bridge operable to receive the bus operation and reserve a data buffer for the exclusive use of said bus master in response to said bus operation, said second bridge operable to forward said reservation command to the first bridge if said forwarding bit is set.

30. The computer system of claim 29, wherein said bus command includes a starting address and wherein said bus command and said reservation command together define an address range.

31. The computer system of claim 30, wherein said reservation command includes a count.

32. The computer system of claim 29, wherein if said bus command is a read operation
said second bus bridge fetches data according to the read operation and prefetches data according to said address range, the data being provided to said bus master when requested.

33. The computer system of claim 32, wherein when said bus master reads at a last address of the address range, the second bus bridge cancels the data buffer reservation.

34. The computer system of claim 29, wherein said reservation command is received after a frame signal is asserted.

35. The computer system of claim 29, wherein said second bus bridge and said bus master are coupled to a peripheral component interconnect (PCI) bus.

36. The computer system of claim 29, wherein said bus operation includes a memory read operation for writing to said mass storage system.

3. Detailed Description of Invention

Field of invention:

The invention relates to a method for translating operations from one bus to another bus and more particularly to a buffer reservation mechanism for optimally transferring data between buses.

Prior art:

Personal computers are constantly changing as new technologies evolve and are incorporated into the computer. Performance improvements in the microprocessor and memory have resulted in computers so powerful that they are now capable of performing tasks that before could only be performed by large mainframe computers. However, to fully replace a mainframe computer, the computer must have significant memory and storage capacity supported by a hearty I/O (input/output) subsystem.

Several standardized I/O buses are available to the system designer including: ISA (Industry Standard Architecture); EISA (Extended Industry Standard Architecture); and PCI (Peripheral Component Interface). Today's computers are typically designed with some combination of the three. For moving data between the buses, a bridge device is typically provided.

The bridge device connects to both buses for transferring data between the buses and translating the bus control signals. The buses can be different, or especially in the case of

PCI, the bridge can simply provide an electrical extension to the same logical bus. This electrical separation makes it possible to meet the PCI bus requirement of limiting the number of physical devices on one bus segment, while at the same time not limiting the total number of PCI agents. In PCI bus vernacular, an agent is the term denoting the class of devices connecting to the bus, including master and slave devices.

These buses all support a scheme called bus mastering wherein a device or agent, usually other than the processor, may request an arbiter for control of the bus. If the arbiter grants the agent control, the agent becomes a bus master. The bus master directs its operations to another agent called a slave. The bus master may then perform operations without processor intervention much more efficiently than if the processor were involved. Many times a master on an origination bus will communicate with a slave on a destination bus. If the buses are coupled by a bridge, the performance of the bridge greatly effects the performance of the communication between the master and slave. Thus, it is desirable to optimize this pathway.

In the case of the PCI bus, one method of improving performance is to permit delayed transactions or read posting. More details on the PCI bus and on Delayed Transactions are found in the PCI Local Bus Specification version 2.1 which can be obtained from the PCI Special Interest Group, Hillsboro, Oregon. The PCI Local Bus Specification and its related documentation are hereby incorporated by reference. Delayed transactions permit the bus to be used while a slow device is preparing data in response to a request. Thus, instead of the slow device applying wait states to the bus, the bus may be used for other requests. For a bridge, the destination bus means the interface that was not acting as the target of the original request. A delayed transaction progresses to completion in three phases: the request by the master; completion of the request by the target; and the completion of the transaction by the master.

In the first phase, the master generates a transaction on the bus, the target decodes the access, latches the information required to complete the access and terminates the request with a retry-termination. Since the master cannot distinguish between a target which is completing the transaction using delayed transaction termination and a target which simply cannot complete the transaction at the current time, it must reissue the request. During the second phase, the target independently completes the request on the destination bus using the latched information from the delayed request. If the delayed request is a read, the target obtains the requested data and completion status. If the delayed request is a write, the target

delivers the write data and obtains the completion status. During the third phase, the master successfully re-arbitrates for the bus and reissues the original request. The target decodes the request and provides the master with the completion status (and data if a read request). All bus commands that must complete on the destination bus before completing on the originating bus may be completed as a delayed transaction. These include interrupt acknowledge, I/O read, I/O write, configuration read, configuration write, memory read, memory read line and memory read multiple commands. Memory write and memory write and invalidate commands can complete on the originating bus before completing on the destination bus. These commands are not completed using delayed transactions termination and are normally posted.

One such bridge device is the Intel PCI to EISA bridge chip set. The 82375EB/SB PCI-EISA Bridge and the 82374EB/SB EISA system component work in tandem to provide an EISA I/O interface for computers having a PCI bus. The chip set can be either a master or slave on both the PCI and EISA buses. For PCI to EISA data transfers, four 32-bit posted write buffers are provided to enhance single cycle PCI bus transactions. For EISA to PCI data transfers, four 16-byte line buffers are included to support EISA bursting.

In order to use both buses efficiently, most bridges implement some amount of data buffering within the bridge itself. This allows the bridge to de-couple the buses from each other and let each bus run at its maximum speed without being slowed down by the other. There are generally two types of buffers that may be implemented in a bridge: write posting buffers and read prefetch or read ahead buffers. Both types can be implemented on either bus.

Write posting buffers accept write data from one bus and acknowledge reception to that bus. This frees the bus to perform other transactions. The bridge temporarily stores, or posts, the write data until it can be written to the other bus. Read prefetch buffers take the address from a single read access and read additional data speculating that it will also be needed. The bridge then holds that data in a buffer until it is either unusable or it is used by a read access.

In the Intel chip set, the buffer permits the bridge to receive short bursts of data at peak data transfer rates. For example, if a EISA device requests data from memory on the PCI bus, the bridge can burst four 32-bit data words from the PCI memory into its buffer and then release the PCI bus to other PCI requestors while the EISA device reads the buffers. Therefore, the PCI bus is not held up by the EISA device. While the buffer is filled, the

EISA device is notified to read the data from the buffer. The EISA device may then read the data from the buffer at its burst transfer speed. Thereafter, the bridge attempts to keep up with the EISA bus transfer rate by performing short burst of data across the PCI bus to keep the buffer filled. However, if another EISA device requests data from memory, the buffer must be flushed and filled again.

Buffers also introduce problems with data consistency. While data is buffered in the bridge, a bus agent and the processor may have different ideas about what is really in memory. When a bus master issues a read request through the bridge to a target memory range on the other side of the bridge, the bridge must balance between two conflicting goals: performance and data integrity. If the memory operation and transfer were optimized for performance, the bridge would preferably hold large blocks of prefetched data. However, if data is prefetched in large blocks but unused, the bridge must guard against providing stale data on a subsequent read request by a master. Additionally, prefetching large blocks of unused data hurts performance by wasting bus bandwidth. Thus, it is desirable to find a solution which both meets performance concerns and also guarantees data integrity.

Buses can also operate at different frequencies. The bridge must therefore synchronize the signals as they are translated from one bus at one frequency to another bus at another frequency. If a continuous flow of data is to be achieved, the synchronization must not be at the expense of performance.

Means to solve problems:

A bus bridge for connecting two buses includes a plurality of data buffers. The data buffers are used for storing read data, storing prefetched or read-ahead data or for storing write posted data. According to the invention, the use of these data buffers reduces the number of times a device will access the buses, thereby reducing bus latency.

The bridge further includes a mechanism for reserving a data buffer. Two alternative mechanisms are disclosed whereby a bus master or agent coupled to the bridge may request that the bridge reserve a buffer for the bus master's intended data read transfer. If a hierarchy of bridges is present, the reservation request is passed on to upstream bridges supporting the buffer reservation mechanism.

The reservation request received by the bridge includes the addressing information necessary for the bridge to read-ahead or prefetch data as the bus master removes data from the reserved data buffer. When a reservation is placed, three 128-byte data buffers are reserved for the bus master. Data is prefetched into these data buffers in order and on a

rotating basis as read by the master. When a data buffer is first accessed by the master, a next data buffer is prefetched. When the complete address range has been read and the data has been removed by the bus master, the reservation is canceled. If the buffer has not been accessed by a bus master within a programmable period of time, the reservation is also canceled. A reservation algorithm is disclosed for managing the reserved buffers.

The bridge further includes two alternative algorithms for managing the buffers if the reservation mechanism is not available to an agent. With the default algorithm, when a new bus master attempts a read without first reserving a buffer, the bridge allocates a buffer, comprising three 128-byte data buffers, to this master if a buffer is available, otherwise the master is forced to retry the read operation. Once the buffer is allocated, the bridge reads the requested data into a first of the 128-byte data buffers and subsequently provides the data to the bus master. If the master burst reads past the end of the first data buffer and the next read operation from the master starts where the last read ended, then prefetching is enabled for this master. On the next sequential read the bridge will again read the requested data into a second data buffer and subsequently provide it to the master. At this point prefetching is already enabled. Whenever the master again burst reads past the end of one of the buffers, the bridge automatically prefetches data for the next 128-byte buffer. The three 128-byte data buffers are filled in order and on a rotating basis. If the bus master requests a non-sequential read or if the CPU performs a write transaction which might affect the master, the buffer is flushed.

A conservative algorithm is similar to the default algorithm, but the bridge does no prefetching. If the bus master does not read all of the data, the remaining data is flushed when the master disconnects from the bus.

Embodiment:

Referring now to Figure 1A, the bridge device 106 is illustrated according to the preferred embodiment. The bridge device 106 couples a primary bus 110 to a secondary bus 112 of a computer system 100. The computer system 100 further includes a CPU 102 and memory 104 coupled to the primary bus 110. In the preferred embodiment, the CPU 102 includes a processor (not shown) such as a 486, Pentium7 or 586 class processor by Intel or their equivalents and standard computer peripherals and support logic (not shown) as is common in IBM compatible PCs. It is understood that other processors and peripherals could of course be utilized.

A bus master 108, such as a hard disk controller, video graphics controller or network interface controller, and bus slave 114, such as a serial port or parallel port are further connected to the secondary bus 112. For purposes of the principles disclosed in this invention, it is not necessary to identify the specific functionality of the bus master 108 and slave 114, therefore, they are simply referred to herein as generic bus masters and slaves.

The primary and secondary buses 110 and 112 are standard input/output buses such as the peripheral component interconnect (PCI) bus. It is understood that the principles herein disclosed may be applied to other buses as well, such as the industry standard architecture (ISA), extended industry standard architecture (EISA) and microchannel architecture (MCA) bus. Further, as can be appreciated, the same bridge design principles may also be applied to couple two dissimilar buses, such as a PCI bus and an EISA bus. In the preferred embodiment, the primary bus is a local PCI bus with limited electrical fanout and the bridge 106 is used to extend the functionality of the primary bus 110 by way of the secondary PCI bus 112.

In an alternative to Figure 1A, shown in Figure 1B, the secondary bus 112 is in a remote location from the computer system 100. The bridge 106 is further comprised of an upstream part 106a and a downstream part 106b. A cable 118 couples the upstream and downstream bridge instances together. The downstream part 106b is coupled to the secondary bus 112 for communicating with a plurality of bus masters 108a-d. Each bus master is further coupled to a plurality of hard disk drives 116a-d for providing mass storage to the computer system 100. In both Figures 1A and 1B, the bridge 106 provides a common interface for devices on the secondary bus, such as the bus master 108, to communicate with devices on the primary bus, such as the memory 104. The bus master 108 may also communicate with other secondary bus devices, such as the bus slave 114, without requiring any bandwidth from the primary bus 110. Thus, in addition to being a translation unit and bus extender, the bridge 106 also functions as a isolation buffer.

When large blocks of data are transferred through the bridge, such as when the bus master 108 reads or writes to the memory 104, the bridge provides multiple buffers for read prefetching and write posting to enhance the overall system performance. This helps eliminate the typical requirement that both buses be communicating with the bridge 106 before a transfer can take place. The read prefetching is managed according to several alternative algorithms. One algorithm includes a mechanism in which a bus master can reserve a read prefetch buffer for very efficient bus transfers. Other adaptive algorithms handle those cases where the bus master does not request a buffer reservation. However, even for these cases, the algorithms attempt to reserve a buffer and prefetch data when certain conditions signify that a read prefetch will have a high probability of success.

Now turning to Figures 2A and 2B, a block diagram of first and second embodiments of the bridge device 106 are illustrated. As can be seen, the primary bus couples to a primary

bus interface 200 and the secondary bus couples to a secondary bus interface 202. The bridge 106 contains a buffer region 204 for receiving incoming data, read prefetching or read-ahead, or for write posting. The buffers are configured according to the type of operation, as described below.

Specifically, in the bridge 106, four 128-byte write posting buffers 218 are provided to enhance the performance of write operations. When the bus master 108 performs a write operation to main memory 104, the data is written or posted into the write buffers 218 and the bus master 108 receives a completion indication from the bridge 106 before the data is actually written to main memory 104. The bridge 106 then subsequently completes the write operation to main memory 104 without further delay. Thus, the time the bus master 108 controls the secondary bus 112 is greatly reduced since it does not wait for the actual completion indication.

The buffer region 204 also includes eight default buffers 210 for enhancing bus master read operations from main memory 104. A default buffer management algorithm manages the default buffers 210 for storing read data and prefetching read data under certain conditions. The default buffer management algorithm balances between holding the prefetch data as long as possible and yet avoids the possibility of retaining stale data. The buffers are filled and emptied using three 3-bit pointers to 32-bit double word chunks of data (hereinafter called DWORDS) in the buffers. A buffer input pointer (BUFF_IN_PNTR) points to where the next Dword will be stored; a buffer output pointer (BUFF_OUT_PNTR) points to where the next Dword will be read from; and a buffer valid pointer (BUFF_VALID_PNTR) indicates which addresses between BUFF_IN_PNTR and BUFF_OUT_PNTR are valid. Only certain addresses may be valid since to prevent bus locked situations a read completed after a posted memory write can be stored in the buffer, but cannot become valid until after the posted memory write completes.

As an alternative to the default buffer management algorithm, a conservative algorithm is provided to manage the default buffers 210. The conservative algorithm never prefetches and generally only reads data that has been requested by a bus master. Thus, stale data is avoided altogether.

The buffer region 204 also includes a pool of eight reservation resources 220 for enhancing bus master read operations from main memory 104. These resources 220 are managed by a reservation algorithm, discussed below. Each reservation resource 206 consists of three 128-byte data buffers and a corresponding reservation information block

208. Each information block 208 includes the registers necessary to control the reservation buffers, which include: a reservation base address register (RSRV_BASE_ADDR) containing a 64-bit starting address of the reservation; a reservation limit address register (RSRV_LIMIT_ADDR) containing a 64-bit address one larger than the last address of the reservation; a discard timer value programmed to one of four values (off, 25ms, 50ms, or 100ms) for automatically canceling inactive reservation resources; a last buffer indicator register (LAST_BUFF) containing a 5 bit pointer to the last line buffer data was read into; a buffer start pointer register (BUFF_START_PNTR) containing a 64-bit starting address of each line buffer; a buffer limit pointer register (BUFF_LIMIT_PNTR) containing a 64-bit address one larger than the last address of the line buffer; a buffer input pointer register (BUFF_IN_PNTR) for holding the address where the next Dword is stored; a buffer output pointer register (BUFF_OUT_PNTR) for holding the address where the next Dword will be read from; a buffer valid pointer register (BUFF_VALID_PNTR) for indicating which addresses between the BUFF_IN_PNTR and BUFF_OUT_PNTR are valid; a buffer state register; and a buffer prefetch flag register.

The default buffers 210 and reservation resources 220 may be allocated to each requesting master or to each process or thread, thereby allowing a single master to have multiple buffers. The buffer region 204 is organized as a FIFO, and of course, may be sized according to the requirements of the master for optimal bus utilization. Since the bridge 106 is preferably transparent to the requestor and its target, it is desirable for the buffer region 204 to also be transparent to the transaction. However, to optimize the operation of the buffer region 204 it is desirable to control the buffers so that data will be prefetched efficiently and the prefetched data will be coherent with data in main memory. Therefore, to optimize the buffers for the pending transactions, it is desirable for the bridge to receive the range of addresses the master intends to read. Thus, once a bridge obtains the address information, the bridge may freely prefetch data up to the end of the address range. Two alternatives are provided for requesting a reservation resource.

In a first alternative reservation mechanism, illustrated in Figure 2A, the secondary interface 202 contains a command decode block 212a. It is noted that while the command decode block 212a is located in the secondary interface 202 because memory read requests are from downstream bus masters, it is contemplated that a command decode block could also be effective in the primary interface 200 for read operations in the opposite direction. The reservation resource is most effective for read operations initiated by a bus master

residing on the secondary bus 112 from main memory 104, thus, the following description is presented in that context.

A reservation resource 220 can be reserved by a bus master 108 when a special bus command is received by the bridge 106. At this point, the reader may also want to refer to Figure 4A. When the bus master 108 desires to reserve a reservation resource 220 for read transactions from the target memory 104, the bus master 108 issues a first read transaction containing a reservation command (RSRV_COMMAND) for indicating a total number of bytes to be read from the target memory 104. The RSRV_COMMAND is a special PCI bus command represented by the command/byte enable (C/BE#) signals on the PCI bus and is presented to the bridge 106 on the first read transaction when a frame (FRAME#) signal is asserted to indicate the beginning of a transaction, as shown at clock 2 of Figure 4A. The reservation command also includes a count on the AD[31:0] lines to indicate the number of bytes desired. From clock 3 forward, the PCI transaction continues as normal: the type of PCI transaction is provided on the C/BE# lines along with a starting address on the address/data (AD) lines; and data is transferred when available (as shown in clocks 5, 7 and 9). Thus, the count (clock 2) and starting address (clock 3) indicate a range of memory to read. It is understood that there are alternative ways to indicate this range for reserving a buffer, such as by indicating a starting address and ending address. These alternatives are contemplated for the reservation command.

For 64 bit addressing, illustrated in Figure 4B, the RSRV_COMMAND precedes the dual address command and actual read command. More details on the dual address command is found in the PCI Specification referred to above.

When the bridge 106 receives the RSVR_COMMAND, the registers in the reservation information block 208 are initialized. The starting address is stored in the RSRV_BASE_ADDR, the ending address is determined by adding the starting address to the total byte count, which is incremented by one and stored in the RSRV_LIMIT_ADDR. One of the three 128-byte buffers is initialized by storing the starting address in the BUFF_START_PNTR and the buffer ending address in the BUFF_LIMIT_PNTR. Also, the buffer state and prefetch flag for each of the three buffers is initialized.

Once the reservation is placed, the bus master 108 may disconnect from the bridge 106 and reconnect later to resume the operation without reissuing the RSVR_COMMAND. Thus, in this alternative, a single PCI bus command indicates the number of bytes to be

transferred. The penalty for this additional information is only one clock cycle, since once the buffer is reserved, the data transactions proceed as normal.

The target does not receive the RSRV_COMMAND. However, intermediate bridges may receive the RSRV_COMMAND from downstream bridges. If a hierarchy of bridges is present, as illustrated in Figures ^{5A}5A and ^{5B}5B, each bridge supporting the RSRV_COMMAND may pass the reservation upstream. In the configuration of each bridge, a configuration bit 214 (Figure ^{3A}3A) is cleared if the upstream bridge does not support buffer reservations, otherwise the bit is set by configuration software during computer initialization operations and the reservation will be passed upstream. Figure ^{5A}5A shows an example hierarchy of bridges where a bridge 300 is a non-reservation bridge sandwiched between two bridges 106a and 106b which support buffer reservations. In this example, since the bridge 300 does not support buffer reservations, the configuration bit in bridge 106b would be cleared so that the RSRV_COMMAND would not be forwarded to the non-reservation bridge 300. Thus, even though bridge 106a supports buffer reservations, this feature is not utilized in bridge 106a.

Figure ^{5B}5B illustrates an example hierarchy of bridges in which a top-most bridge 300 does not support buffer reservations, but a lower two bridges 106a and 106b do support it. In this example, since the bridge 300 does not support buffer reservations, the configuration bit in the bridge closest to the non-reservation bridge, bridge 106a, would be cleared so that the RSRV_COMMAND would not be forwarded to the non-reservation bridge 300. The configuration bit 214 in bridge 106b is set so that the RSRV_COMMAND received by bridge 106b is forwarded to bridge 106a. Thus, a non-reservation bridge acts to cut off the benefits of buffer reservations when implemented with the RSRV_COMMAND.

Referring now to a second alternative reservation mechanism, illustrated in Figure ^{4A}4A, the bridge 106 contains an address decode block 212b for decoding operations to a special address. In the memory address map of the computer system 100, this special memory address is allocated for this reservation mechanism. Write operations to this special address are interpreted as reservation operations by the bridge 106. When the bus master 108 desires a buffer reservation for a read operation, it first executes a 2-3 32-bit (Dword) burst write to the special address. The bridge 106 intercepts the write operation and stores the command in a reservation mailbox 216. The first Dword is the total number of bytes the bus master 108 intends to read. The second Dword is the least significant 32 bits of the starting address. If 64-bit addressing is utilized, then a third Dword is written containing the

most significant 32 bits of the starting memory address. If only two Dwords are written, the bridge assumes the most significant 32 bits are zero. The bridge 106 intercepts and claims the write operation to the mailbox register 216 by asserting a device select (DEVSEL#) signal since the address is not assigned to any devices on the secondary bus 112. When the bridge 106 receives the transaction, it reserves a buffer if one is available, and stores the address, length and other information in the mailbox 216 for forwarding to other intermediate bridges, if necessary. When the reservation is placed, the reservation information block is also initialized as described above. Thus, in this alternative, a write to a special address indicates the starting address and number of bytes to be transferred. This information is utilized by the bridge 106 to prefetch data for the requesting master 108.

If a hierarchy of bridges is present as illustrated in Figures ^{5A} 5A and ^{5B} 5B, each bridge supporting buffer reservations may pass the reservation information upstream. In the configuration of each bridge, the topmost bridge in the hierarchy is disabled from forwarding the reservation. Figure ^{5A} 5A illustrates a hierarchy of bridges where a bridge 300 is a non-reservation bridge sandwiched between two bridges 106a and 106b which support buffer reservations. In this example, since write transactions to global memory address are used to pass the reservation information between the bridges, the configuration bit and bridge 106b would be set so that the reservation would be forwarded to the non-reservation bridge 300 and passed to the reservation bridge 106b. The non-reservation bridge 300 does not recognize the difference between a normal memory transaction and a reservation write transaction. If the bridge 106a was coupled to the memory 104 on its primary side, the configuration bit and bridge 106a would be set to disable forwarding of the reservation information to memory. Thus, each bridge capable of supporting buffer reservations receives the reservation information.

Figure ^{5B} 5B illustrates a hierarchy of bridges in which a topmost bridge 300 does not support buffer reservations, but a lower two bridges 106a and 106b do support it. In this example, since the bridge 300 is the topmost bridge and because it does not support buffer reservations, the configuration bit in bridge 106a would be cleared to disable forwarding of the reservation information. The reservation is passed between the bridges 106a and 106b. Thus, if a non-reservation bridge is located at the top level, the reservation bridge below the non-reservation bridge is disabled from forwarding the reservation. One advantage of this method over the first alternative is that script-driven bus masters can readily implement the memory write transaction with little modification.

Both alternatives allow the bridge 106 to prefetch data in longer, more efficient transactions. By utilizing the buffer reservation, unused prefetched data is avoided since only the data requested is prefetched. Coherency problems between the memory 104 and data stored in the buffers 206 and 210 are also avoided since only requested data is prefetched and it can be assumed that no bus master would use this data until the requestor completes the transfer. Thus, data prefetched according to the preferred embodiment can be held and considered valid until the transfer completes.

Read operations may be performed by the bridge 106 with a technique known as delayed transactions. More details on delayed transactions are found in the PCI Specification incorporated above and in the background section of this specification. Briefly though, delayed transactions are a method of termination and retry used by PCI bridges and PCI I/O controllers which cannot complete a data transaction within a certain specified time. Instead of a first requesting bus master waiting for a slow I/O controller or bridge to respond, the I/O controller or bridge may force the first requesting bus master off the bus so that other bus masters are allowed to use the bus bandwidth that would normally be wasted holding the first bus master in wait states. A delayed transaction normally progresses to completion in three phases: a request by a bus master; completion of the request by a target; and a completion of the transaction by the bus master.

In the nomenclature of the PCI bus, a delayed read request is a transaction that must complete on the destination bus before completing on the originating bus and can be an I/O read, configuration read, memory read, memory read line, or memory read multiple command. Referring to Figure 4A, if the bus master 108 requests data from memory 104, the secondary bus 112 is the originating bus and the primary bus 110 is the destination bus. Once a request has been attempted on the destination bus, it must continue to be repeated until it completes on the destination bus. Until that time, the delayed read request is only a request and may be discarded at any time to prevent bus deadlock or improved performance since the bus master 108 must repeat the request later. A delayed read completion is a read transaction which has completed on the destination bus and is now moving toward the originating bus to complete. The delayed read completion contains the data requested by the bus master 108 and the status of the target (memory 104). Once the delayed read request is executed on the destination bus it becomes a delayed read completion.

An example of a simple delayed transaction without the buffer reservation mechanisms is helpful to illustrate its effect on the buffers region 204. It is noted that

delayed transactions may be used with or without the buffer reservation mechanisms herein disclosed. Also, the buffer reservation mechanisms may be used with or without delayed transactions, but it is assumed they will be used together for maximum performance.

Figure 4A is used as an example. During a first phase, the bus master 108 generates a memory read transaction on the secondary bus 112. The bridge 106 decodes the transaction and, recognizing that the main memory is on the primary bus, claims the transaction by asserting the DEVSEL# signal, latches the transaction information required to complete the access, and terminates the request using conventional PCI retry techniques. The latched request information is referred to as a delayed request, or more specifically in this example, a delayed read request. According to PCI convention, the bus master 108 must repeat the transaction. Meanwhile, using the latched transaction information, the bridge 106 generates a memory read transaction on the primary bus 110. If the target (main memory 104) is not ready, it can also force another delayed transaction by issuing a retry to the bridge 106. However, assuming the memory 104 does not issue a retry, the memory 104 processes the request and the bridge 106 receives the completion status and the requested data into one of its buffers 204. If the delayed request was a write operation, the bridge 106 would obtain the completion status from the main memory 104. The result of completing the delayed request on the destination bus (primary bus 110) produces a delayed completion, which consists of the latched information of the delayed request and the completion status and data. The bridge 106 stores the data and completion status in the buffer region 204 until the bus master 108 repeats the initial request.

During the third phase, the bus master successfully rearbiterates for the bus and reissues the original memory read request. The bridge 106 decodes the request, claims the operation and provides the data to the bus master 108. At this point, the delayed completion is retired and the transaction has completed. The status returned to the bus master 108 is exactly the same as the bridge 106 obtained from the target (memory 104) when it executed the read transaction, i.e., master-abort, target-abort, parity error, normal, or disconnect.

By using the delayed transaction technique in combination with the buffer reservation mechanisms herein disclosed, read operations between bus masters and bus slaves can be performed more efficiently than before.

Managing the buffer region 204 is desirable to obtaining efficient data transfers. The basic component of the buffer region 204 is the 128-byte buffer, referred to generically as a buffer 500 in Figure 5A. Read prefetch data and posted write data may both be stored in

the bridge 106 (in different buffers), but the read prefetch data cannot become valid until after the posted memory write completes. As data is written into the buffer 500, a `BUFF_IN_PNTR` is incremented to indicate the next address. Data is emptied from the buffer 500 beginning at the address indicated by `BUFF_OUT_PNTR`. The `BUFF_VALID_PNTR` indicates which region of the buffer 500 is valid. Addresses between `BUFF_OUT_PNTR` and `BUFF_VALID_PNTR` are valid addresses. When a delayed read request is started, the `BUFF_IN_PNTR` and `BUFF_OUT_PNTR` are initialized to bits 4-2 of the PCI transaction address. As each Dword of data arrives, it is stored at the address indicated by the `BUFF_IN_PNTR` and the pointer is then incremented. If there are no preceding unexecuted posted memory write addressing targets on the destination bus, the `BUFF_VALID` pointer can be incremented when `BUFF_IN_PNTR` is incremented. If there is a preceding posted memory write, then `BUFF_VALID_PNTR` is updated after the write completes. As the bus master 108 takes data from the buffer, `BUFF_OUT_PNTR` is incremented after each Dword is removed.

Figure 9B illustrates how the reservation resource buffer 206 is constructed. The reservation buffer 206 has a beginning address, as indicated by `RSRV_BASE_ADDR` and ending address, as indicated by `RSRV_LIMIT_ADDR`. The reservation buffer 206 is comprised of three 128-byte data buffers (of the type shown in Figure 9A). As a read transaction progresses, these three buffers are reused so as to simulate one long buffer of length `RSRV_LIMIT_ADDR` minus `RSRV_BASE_ADDR`. Each individual buffer has a starting address, as indicated by `BUFF_START`, and an ending address as indicated by `BUFF_LIMIT`.

The state of the buffer is indicated by two variables: `BUFF_STATE` and `PREFETCH_FLAG`. The `PREFETCH_FLAG` is used to distinguish between a prefetch operation and an actual read request. The state of the data stored in the buffers is illustrated in figure 10A. The state of the buffer is designated as `EMPTY` after a system reset, a buffer timeout or buffer flush. When the buffer is designated `EMPTY` it is available for allocation. If the buffer is servicing a read request or is performing a read prefetch, the buffer is designated as `REQUEST`. Once some data begins to arrive into the buffer, the state transitions from `REQUEST` to `PART_COMP`. Once all the data has been received, the state of the buffer transitions from `PART_COMP` to `COMPLETE`. If the buffer's state is `REQUEST`, `PART_COMP` or `COMPLETE` and the buffer is flushed, the state returns to

EMPTY. If the buffer is designated as COMPLETE and the bus master reads all of the data from the buffer, the state is changed from COMPLETE to EMPTY.

For the reservation algorithm, the buffer is only flushed when the last DWORD of data is removed from the buffer. However, flushing of prefetch data is done from time to time with the default or conservative algorithms to guarantee that a bus master is never given stale data. For this reason, it is important to track the state of the buffers.

Prefetch data refers to the remnant of a buffer which was filled with a real delayed read request from a master, but when the master reconnected only part of the data was taken before the master disconnected. A master's prefetch data is flushed under the following conditions: (1) the master attempts a read from an address which is not sequential from that master's last read; (2) the CPU 102 does a write transaction which might affect this master such as an I/O write to a secondary bus target, a memory write to the general secondary memory range register, or a memory write to the slot-specific memory range register for this master; and (3) a secondary bus master does a write transaction which might affect this master such as an I/O write to a secondary bus target, a memory write to the general secondary memory range register, a memory write to the slot-specific memory range register for this master, or a CPU 102 write to any configuration register.

The states for the PREFETCH_FLAG are described with reference to Figure 11. When the BUFF_STATE is designated REQUEST, PART_COMP or COMPLETE, the PREFETCH_FLAG indicates whether the read is a prefetch or an actual REQUEST. When the PREFETCH_FLAG is set, a prefetch operation is in progress, while when the PREFETCH_FLAG is cleared a real REQUEST is in progress. The transitions from the set and clear states are indicated in the following equations:

- A = (prefetch) OR
(after master read AND BUFF_STATE=(PART_COMP OR COMPLETE));
- B = (flush AND BUFF_STATE=(REQUEST, PART_COMP OR COMPLETE));
- C = (master read transaction AND delayed read request) OR
(prefetch complete AND BUFF_STATE=EMPTY) OR
(flush AND BUFF_STATE=COMPLETE) OR
(real delayed read request);
- D = (master read transaction AND delayed read request) OR
(prefetch complete AND BUFF_STATE=EMPTY);

A buffer 204 is only available for a new REQUEST if it is EMPTY or it contains prefetch data and there is not already a prefetch in progress for that buffer. When a bus master 108 attempts a read and the bridge 106 initiates a delayed read request, BUFF_STATE is changed to REQUEST, and PREFETCH_FLAG is cleared. If the read operation is initiated by a bridge 106 prefetch, then BUFF_STATE is changed to REQUEST, and PREFETCH_FLAG is set.

When the first data arrives BUFF_STATE is changed to PART_COMP. If the master repeats the read transaction and the state is PART_COMP or COMPLETE, read data is provided to the bus master 108. If the bus master 108 does not take all of the data before disconnecting from the bridge 106, PREFETCH_FLAG is set to indicate that the remainder of the data in the buffer 204 is now considered to be prefetched data, but the state remains unchanged. If the bus master 108 takes the last data, i.e., (BUFF_OUT_PNTR > = BUFF_IN_PNTR > after the read) when the state is COMPLETE, the state is changed to EMPTY.

If a flush event occurs for the buffer 204 and its PREFETCH_FLAG is set, and the state is designated REQUEST, PART_COMP or EMPTY, a prefetch is currently in progress, so the state is changed to EMPTY, but the PREFETCH_FLAG remains set until the pending prefetch is complete. The buffer 204 cannot be reused for another prefetch until the PREFETCH_FLAG is cleared. If a flush event occurs for the buffer 204 and its PREFETCH_FLAG is set, and the state is designated COMPLETE, the state is changed to EMPTY and the PREFETCH_FLAG is cleared.

If a real delayed requested needs the buffer 204 while the PREFETCH_FLAG is set, the new request is given the buffer. The PREFETCH_FLAG is cleared (to cause incoming prefetch data to be discarded) and the state is changed to REQUEST.

When prefetch data arrives for the buffer 204 and PREFETCH_FLAG is cleared, the data is discarded (the buffer has been reused for a real request). When prefetch data arrives for a buffer 204 and the state is EMPTY, the data is discarded (the buffer is flushed). When the last prefetch data arrives, the PREFETCH_FLAG is cleared and the buffer 204 is available to be reused.

Three prefetch and flush algorithms are implemented in the bridge 106: default algorithm, conservative algorithm, and reservation algorithm. The default algorithm should be applicable to most bus masters. The conservative algorithm prefetches less and flushes more, but is implemented only in the case where the default algorithm does not work for a

particular bus master. The reservation algorithm is used by more sophisticated bus masters which are capable of initiating the reservation mechanism.

Now referring to Figures ⁽¹²⁾ 7A through ⁽¹⁹⁾ 7H, the prefetch and flush algorithms are described according to the preferred embodiment. It is understood that the bridge 106 performs various processes related to read and writing data, however, for simplicity these steps have been omitted. Further, although the steps relating to write posting are not shown, it is understood that the bridge performs these operations in conjunction with the read prefetch and flush processes described herein.

The prefetch and flush algorithms begin at Figure ⁽¹²⁾ 7A with the receipt of a bus transaction on the secondary bus interface 202. At step 700 (Fig. ⁽¹²⁾ 7A), it is determined whether a configuration command was received. If so, control proceeds to step 710 (Fig. ⁽¹³⁾ 7B) to complete the command, otherwise control continues to step 702. At step 702 it is determined if a write operation has been issued to the bridge 106. If so, control proceeds to step 712 (Fig. ⁽¹⁴⁾ 7C) where it is determined if the write operation affects one of the default buffers that have been allocated. If a bus master 108 is allocated one of the default buffers 210, and a write operation is performed to the bus master 108, such as a configuration write or command, then it is possible that a new process has been initiated, in which case any prefetched data contained in the default buffer 210 is subject to being stale or non-conforming. In such a case, the buffer is said to be affected by the write command. If such is the case, then control proceeds to step 714 where the default buffer 210 corresponding to the slot or master affected by the write operation is flushed. Control then proceeds from step 714 to step 716. If it is determined at step 712 that the write command does not affect any default buffers 210, then control proceeds from step 712 to step 716. At step 716, the remaining write operation processing is handled, including receiving any write data into the write posting buffers and updating the BUFF_VALID_PNTR (at step 718). The BUFF_VALID_PNTR is only updated for the default and conservative algorithms. The routine then terminates, returning to whatever process called it.

Returning to step 702, if it is determined that the operation is not a write operation then control proceeds to step 704 where it is determined if the operation is a reservation command. If so, control proceeds to step 720 (Fig. ⁽¹⁵⁾ 7D) where the reservation command is received by the bridge 106. Control then proceeds to step 722 where it is determined if one of the reservation resources 220 is available for this request. If not so, then control proceeds to step 742 where a retry is issued to the requesting bus master, and the routine ends. If a

reservation resource 220 is available, then control instead proceeds from step 722 to step 724 where the reservation resource registers are initialized with data (i.e., address and count) contained in the reservation command. Control then proceeds to step 728 where three 128 byte buffers are allocated to this request and to step 730 where the status for these buffers is changed to EMPTY. Control then proceeds from step 730 to step 732 where the bridge 106 prefetches three to four lines of data into a first 128-byte buffer of buffer 206. Control then proceeds to step 734 where the last buffer pointer (LAST_BUFF_PNTR) is initialized to indicate which buffer will be written to next.

Control then proceeds to step 736 where it is determined if the forwarding configuration bit is set. If set, then control proceeds to step 738 where the reservation command is forwarded upstream to any upstream bridge. Control then proceeds from step 738 to step 740. If at step 736, the configuration bit is not set, then control proceeds directly to step 740. At step 740, a reservation timer is reset and started. The reservation timer cancels the reservation if after a certain amount of time, such as 25ms, 50ms or 100ms, the reservation resource has not been accessed. It is assumed that the data will be stale after this period of inactivity. The routine then ends.

Returning to step 704, if it was determined that the operation is not a reservation command then control proceeds to step 706 where it is determined if the operation is a read retry from a bus master. If not so, then the operation is deemed to be a first read operation and control proceeds to step 744 (Fig. 7H).

At step 744 it is determined if the requesting bus master already has a default buffer 210 allocated to it or is a new bus master. If it is determined that this is not a new master then control proceeds to step 746 where it is determined if the requested address is equal to the address pointed to by the BUFF_OUT_PNTR. If so, then control proceeds to step 750. If not so, then control proceeds to step 748 where the allocated default buffer 210 is flushed because the master has begun reading from a new, not contiguous address. Control then proceeds to step 750 where the default buffer 210 is reallocated to this master and the buffer registers are initialized as described above. Control then proceeds to step 768.

If at step 744 it was determined that this is a new master, then control proceeds to step 756 where it is determined if any one of the default buffers 210 currently are designated as EMPTY. If so, then control proceeds to step 758 where one of the EMPTY default buffers 210 is allocated to the new master. From step 758, control proceeds to step 768. If not so, then control proceeds to step 760 where it is determined if any default buffers 210

are currently designated as REQUEST or PART_COMP. If so, then control proceeds to step 762 where it is determined if the PREFETCH_FLAG is set and if so, then control proceeds to step 764 where the default buffer 210 is allocated to the new master. Control then proceeds from step 764 to step 768. If at steps 760 and 762, the determinations are negative, then control proceeds to step 766 where the bus master is asked to retry the operation at a later time.

At step 768, control disables prefetching for the newly allocated or re-allocated default buffer 210. Control then proceeds to step 770 (Fig. ⁽¹⁶⁾7E) where the read request information is stored. Control then proceeds to step 772 where the bus master is asked to retry the operation while the bridge 106 gathers the requested data from the target on the destination bus. Control then proceeds to step 774 where it is determined if the read request was a memory read multiple command. If so, then control proceeds to step 776 where the bridge 106 reads three to four lines of data from main memory 104. Control then proceeds from step 776 to step 784. If at step 774 it is determined that the read command was not a memory read multiple, then control proceeds to step 778 where it is determined if read command was a memory read line command. If so, then control proceeds to step 780 where the bridge 106 reads one line of data from main memory 104 and control then proceeds to step 784. If at step 778 it is determined that the command was not a memory read line command, then control proceeds to step 782 where the bridge 106 reads one Dword of data from main memory 104 and then control proceeds to step 784.

At step 784 the bridge 106 receives the data from main memory 104 and at step 786 the bridge 106 stores the received data in the allocated buffer and ends. The routine will later be called when the requesting bus master retries the operation.

Returning now to step 706, if it is determined that the operation is a read retry operation from a bus master then control proceeds to step 708 where it is further determined if this bus master has previously requested a buffer reservation. If so, then control proceeds to step 788 (Fig. ⁽¹⁸⁾7F) where it is determined if the address for the retry read operation is within the limits of the reservation resource 220. If not so, then control proceeds to step 790 where the reservation is canceled and control proceeds to step 744 (Fig. ⁽¹⁷⁾7D) to begin a read operation using the default algorithm. If at step 788 it was determined that the read operation is within the reservation limits then control proceeds to step 792 where it is next determined if the read address is within one of the three buffer limits. If so, then control proceeds to step 794 where the requested Dword of data is provided to the bus master. From

step 794, control proceeds to step 796 where it is determined if this is the first time the bus master has read from this particular 128 byte buffer.

If this is the first time, then control proceeds to step 798 where the bridge 106 prefetches data into the next buffer as indicated by the LAST_BUFF_PNTR. The buffers 206 are filled according to a least recently used algorithm. Once data is read from a buffer for the first time, the least recently used buffer is used for this next data prefetch operation. So at this point, the bus master is reading from one buffer, a next buffer holds prefetched data, and another buffer holds previously prefetched data. So even if the bus master needs to back-up or step back from its current address, the data will still likely be present in the previously prefetched buffer. Control proceeds from step 798 to step 800. If this is not the first time, then control proceeds from step 796 to step 800 where it is determined if the read operation by the bus master is a burst read. If so, then control proceeds back to step 788 to provide additional data. If not so, then control proceeds to the end.

If at step 792 it was determined that the read operation was not within the buffer limits then control proceeds to step 804 where the bridge 106 forces the master to retry the operation while the bridge 106 fetches the requested data. Control then proceeds from step 804 to step 806 where the bridge reads the data from main memory 104, proceeds to step 808 where it receives the data, and proceeds to step 810 where the received data is stored in the buffer 206 and the bridge 106 then awaits for the bus master to retry the read operation.

Returning to step 708, if it was determined that there was not a reservation resource previously requested for this read operation then control proceeds to the default algorithm at step 812 (Fig. ¹⁹7H).

At step 812, the requested data is provided to the bus master from one of the default buffers 210. Control then proceeds to step 814 where it is determined if the read operation is a burst operation. If not so, then processing continues to step 815 where it is determined if the conservative algorithm is selected. If so, then control proceeds to step 817 where the buffer is flushed. If the conservative algorithm is not selected, then control proceeds to the end. If at step 814 a burst operation is in progress, then control proceeds to step 816 to determine if the next requested address is contained within the data buffer. If so, then control proceeds to step 818 where the data is provided to the bus master. Control then proceeds back to step 814 to provide additional data if requested. If at step 816 it is determined that the next requested address is not contained within the buffer, then control proceeds to step 820 where it is determined if prefetching has been enabled for this buffer.

If not so, then control proceeds to step 832 where prefetching is enabled. Control then proceeds to step 770, as described above.

If at step 820 it is determined that prefetching has been enabled, then control proceeds to step 822 where the read transaction information is stored in the bridge 106. Control then proceeds to steps 824 and 826 where the bridge 106 issues a retry to the bus master and the bridge issues a read operation to main memory 104. From this point forward, prefetching is enabled for this bus master. At step 828 the bridge 106 receives the read data and at step 830 the bridge stores the received data into the buffer and waits for the bus master to retry the read operation.

The foregoing disclosure and description of the invention are illustrative and explanatory thereof, and various changes in the size, shape, materials, components, circuit elements, wiring connections and contacts, as well as in the details of the illustrated circuitry and construction and method of operation may be made without departing from the spirit of the invention.

4. Brief Description of Drawings

Figure ⁽¹⁾1A is a block diagram of a computer system incorporating a bridge device according to the preferred embodiment;

Figure ⁽²⁾1B is a more detailed block diagram of a computer system incorporating a bridge device according to the preferred embodiment;

Figure 3 is a block diagram illustrating the functional blocks of the bridge device according to the preferred embodiment;

Figure 4 is a block diagram illustrating the functional blocks of the bridge device according to the preferred embodiment;

Figure ⁽⁵⁾5A is a block diagram illustrating an exemplary hierarchy of buses with a non-reservation bridge coupled between two reservation bridges according to the preferred embodiment;

Figure ⁽⁶⁾5B is a block diagram illustrating an exemplary hierarchy of buses with a non-reservation bridge acting as a host bridge at the host top of the hierarchy;

Figure ⁽⁷⁾6A is a timing diagram illustrating a PCI bus read transaction incorporating a reservation command;

Figure ⁽⁸⁾6B is a timing diagram illustrating a PCI bus read transaction with a 64-bit address incorporating a reservation command;

Figure 9 shows block diagrams of buffers, wherein A is a block diagram of a single 128-byte data buffer according to the preferred embodiment, and B is a block diagram of a reservation buffer according to the preferred embodiment;

Figure ⁽¹⁰⁾6A is a state diagram illustrating buffer states according to the preferred embodiment;

Figure ⁽¹¹⁾6B is a state diagram illustrating buffer prefetch states according to the preferred embodiment; and

Figure 12 is a flow diagram illustrating a prefetch and flush algorithm according to the preferred embodiment;

Figure 13 is a flow diagram which follows "A" of Figure 12, according to the preferred embodiment;

Figure 14 is a flow diagram which follows "B" of Figure 12, according to the preferred embodiment;

Figure 15 is a flow diagram which follows "C" of Figure 12, according to the preferred embodiment;

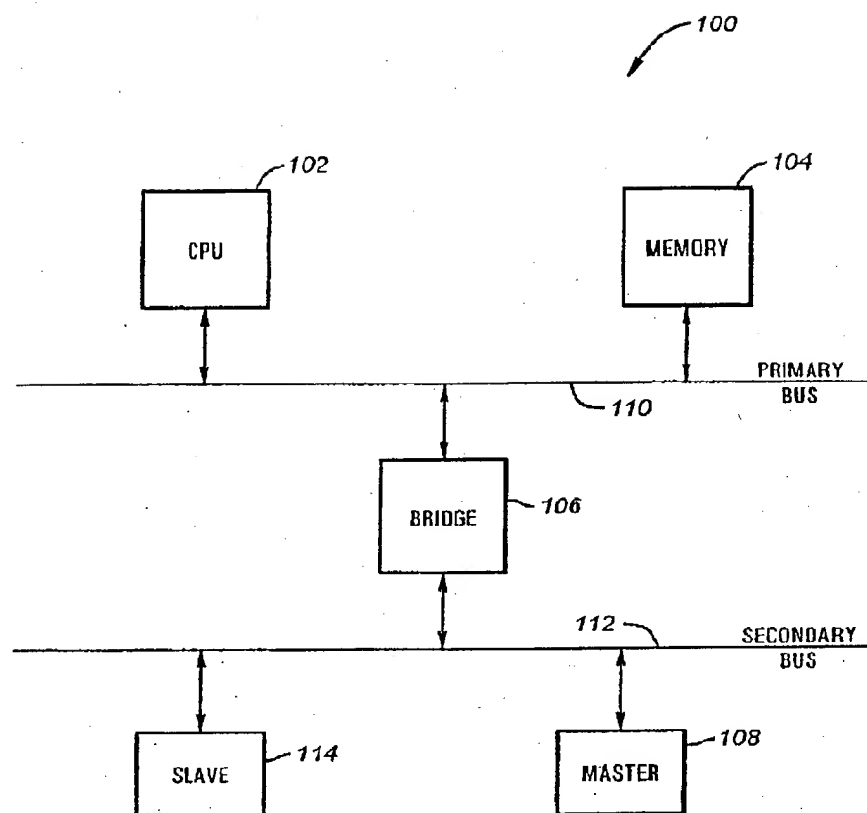
Figure 16 is a flow diagram which follows "G" of Figures 17 and 19, according to the preferred embodiment;

Figure 17 is a flow diagram which follows "D" of Figure 12, according to the preferred embodiment;

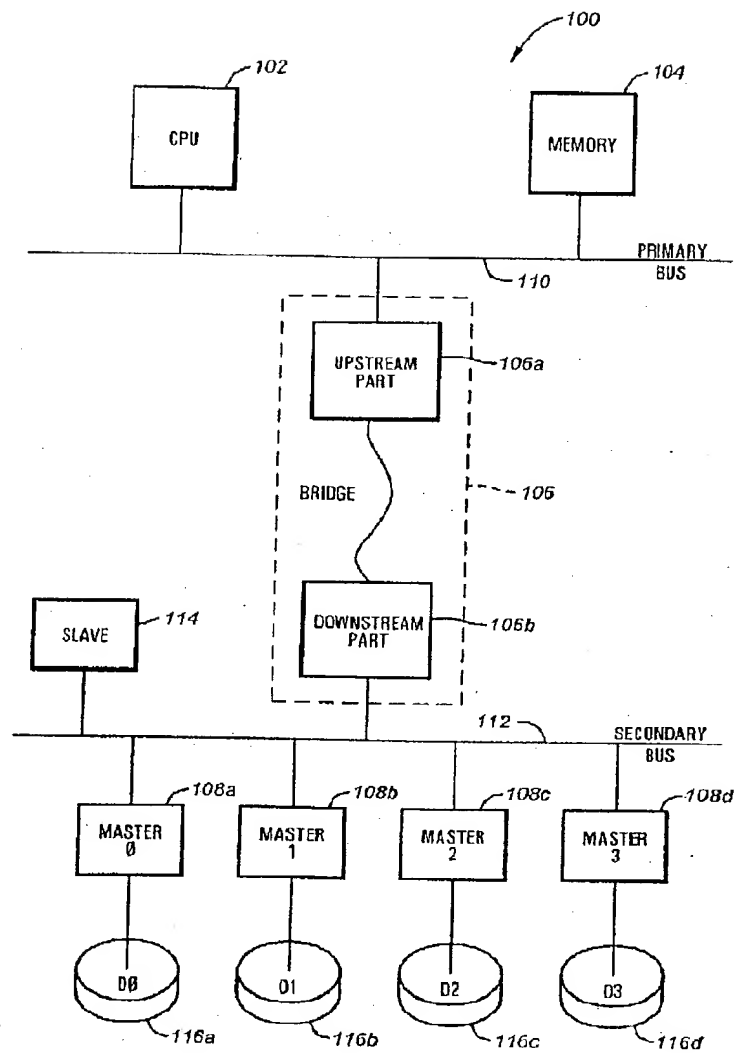
Figure 18 is a flow diagram which follows "E" of Figure 12, according to the preferred embodiment;

Figure 19 is a flow diagram which follows "F" of Figure 12, according to the preferred embodiment;

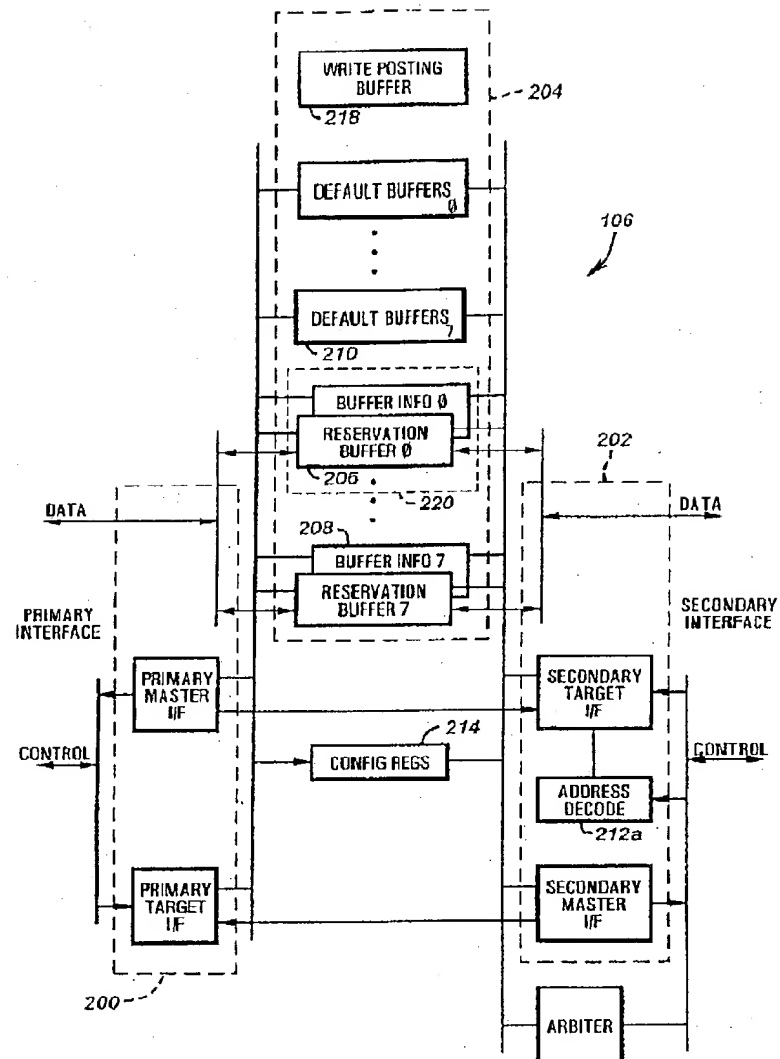
【図1】



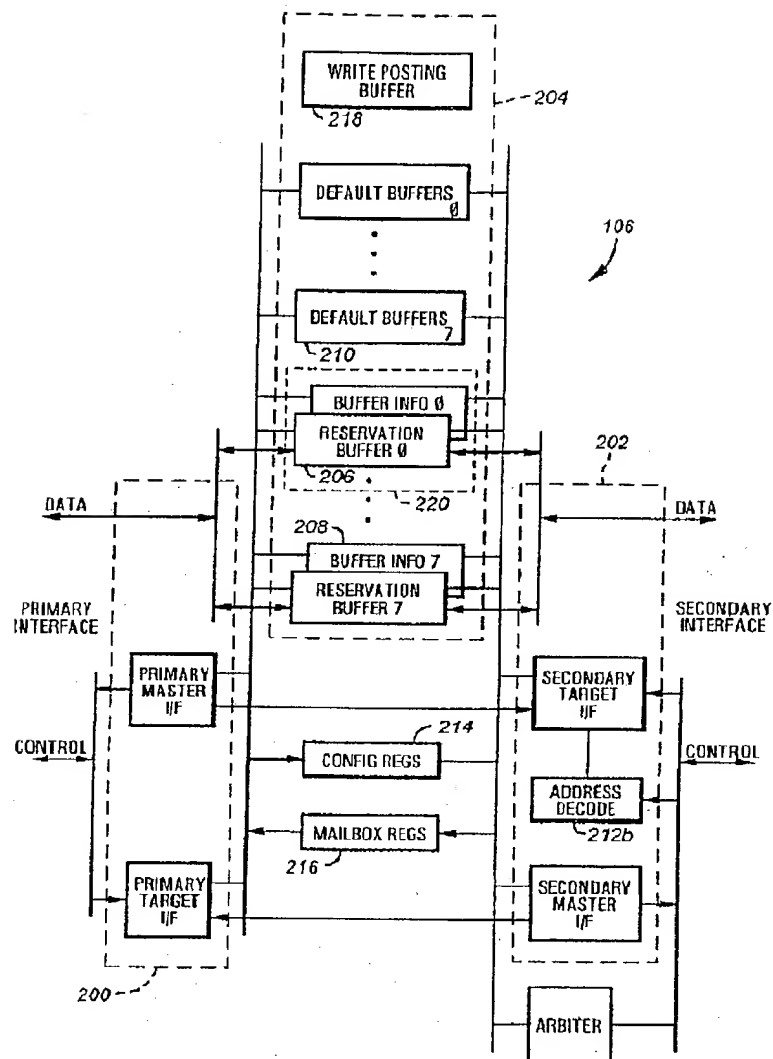
【図2】



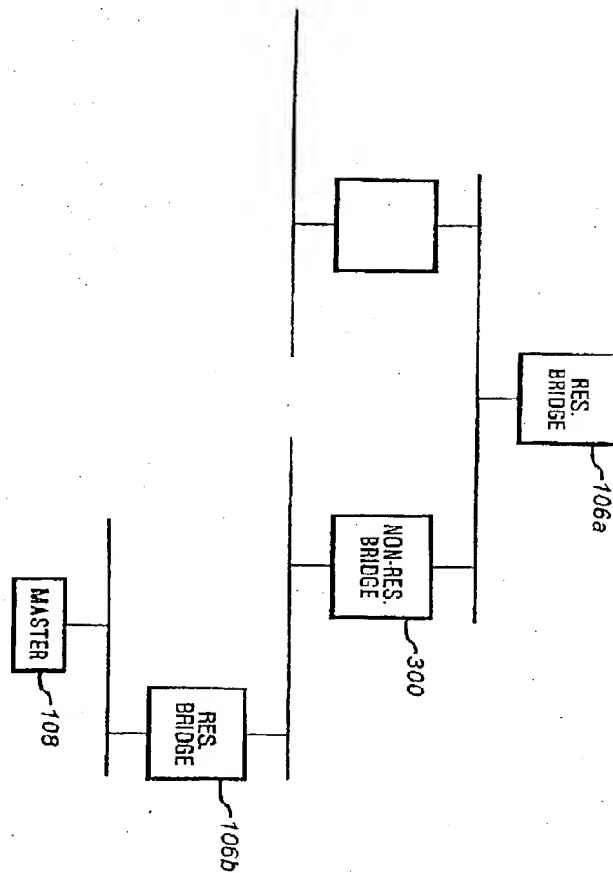
【図3】



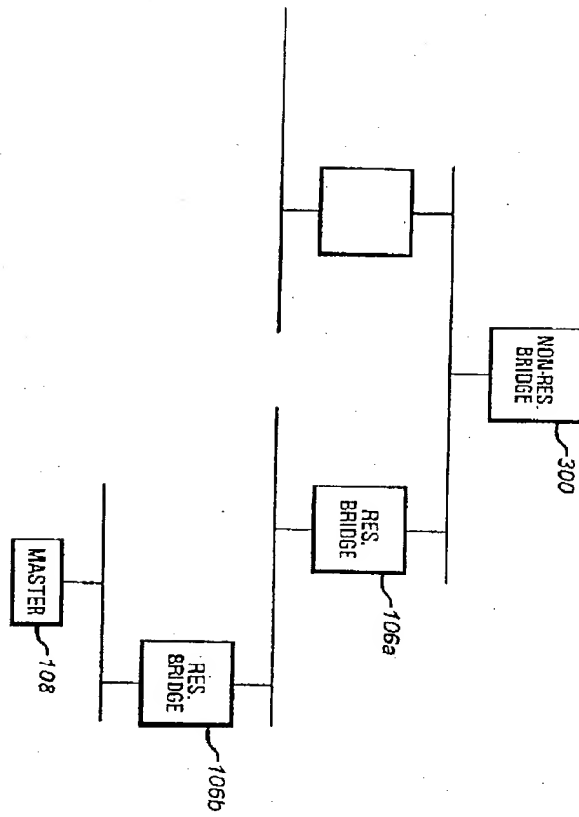
【図4】



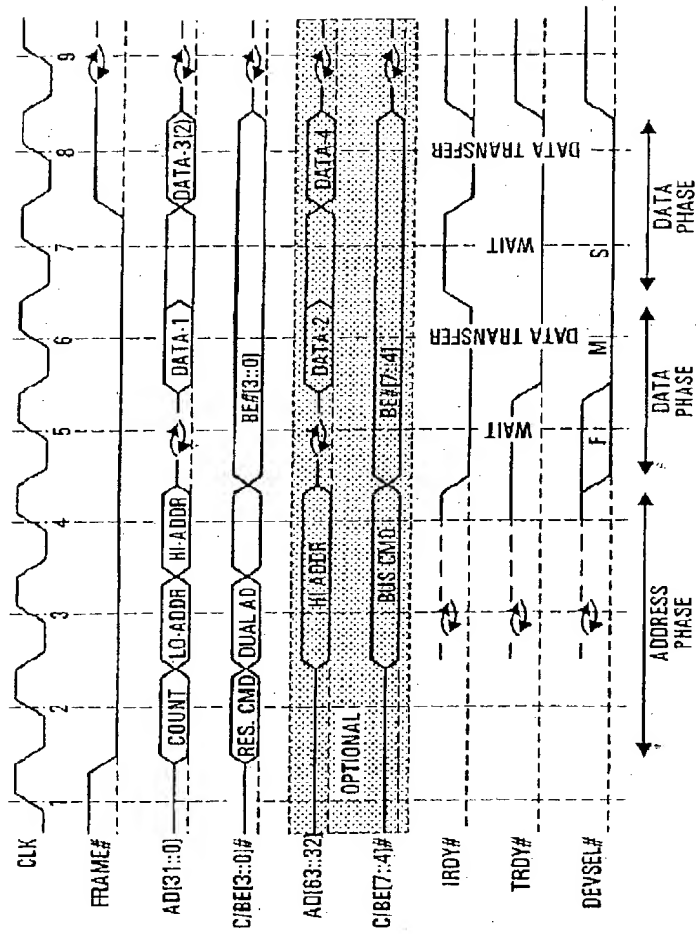
【図5】



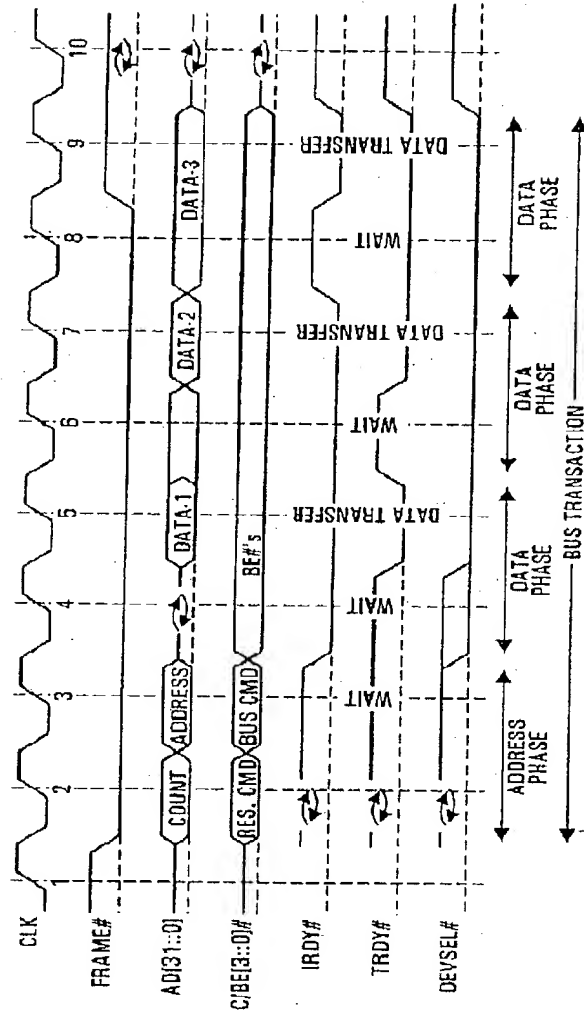
【図6】



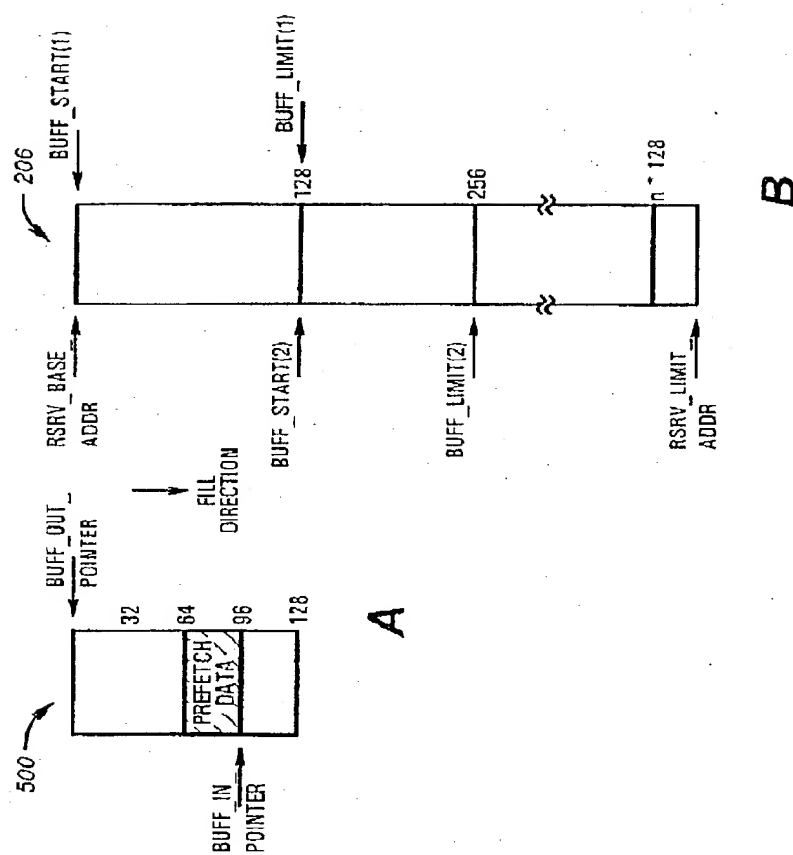
【図7】



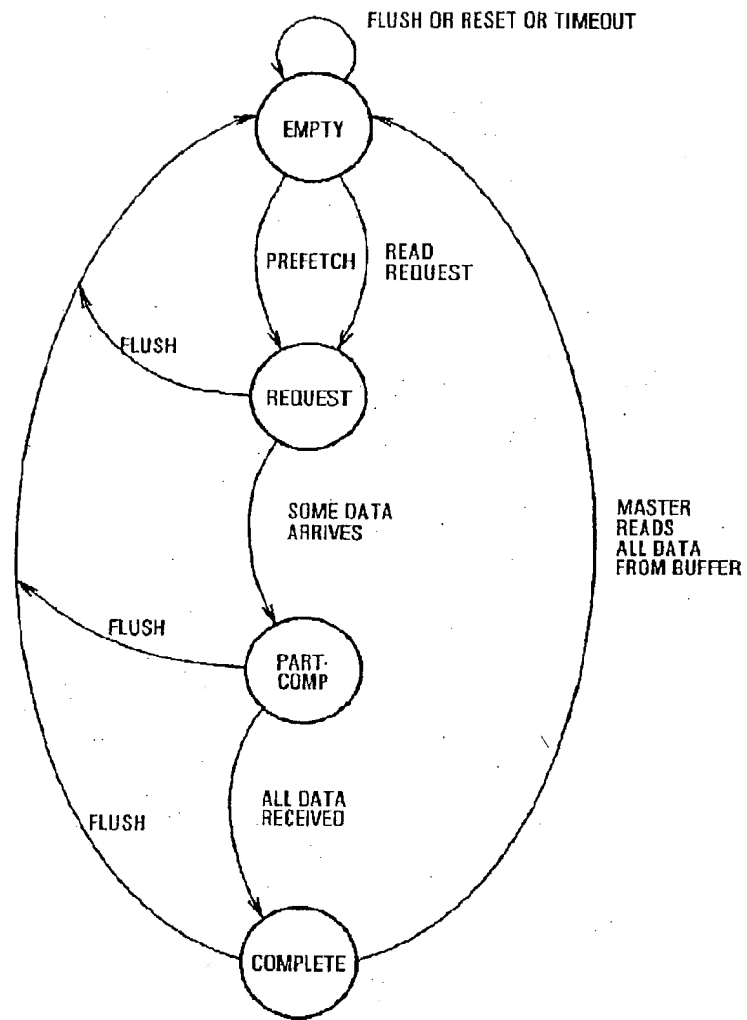
【図8】



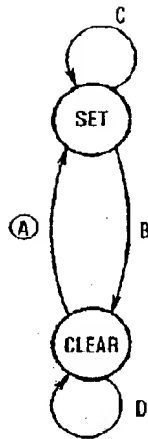
【図9】



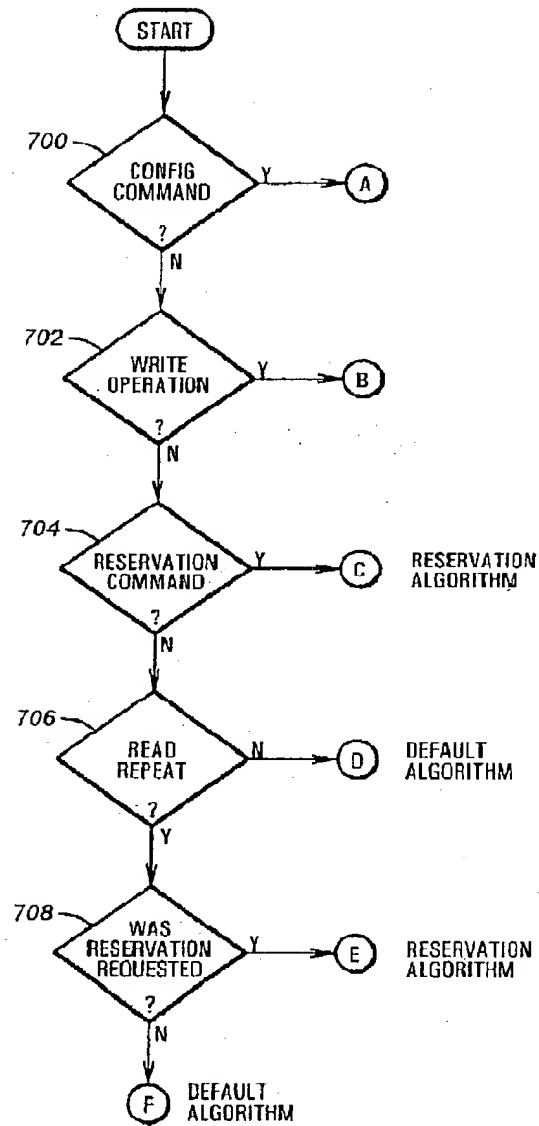
【図10】



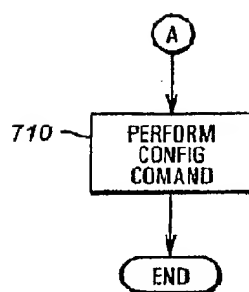
【図11】



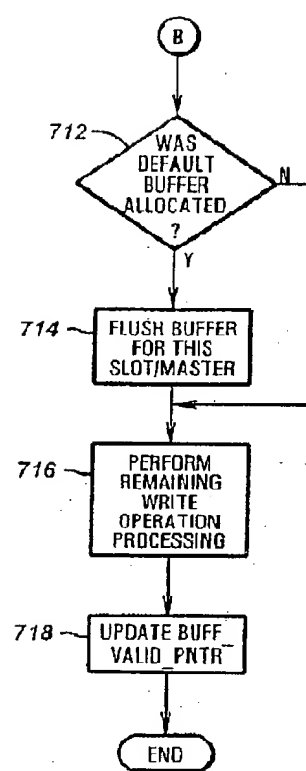
【図12】



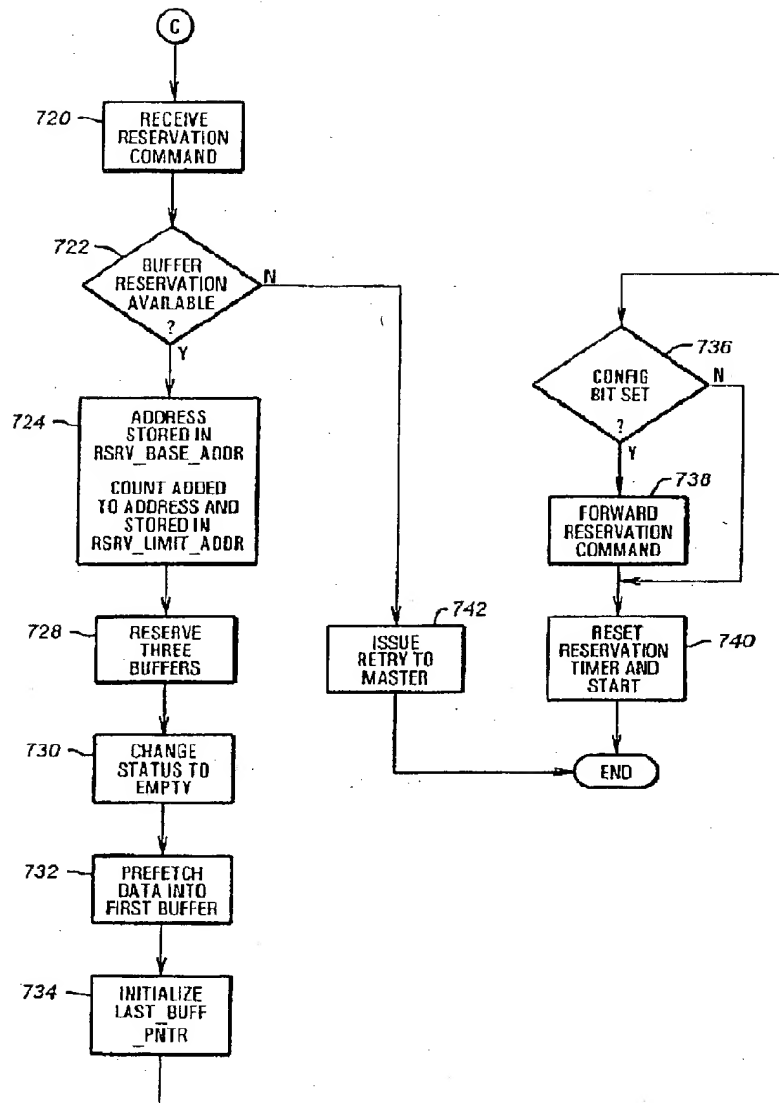
【図13】



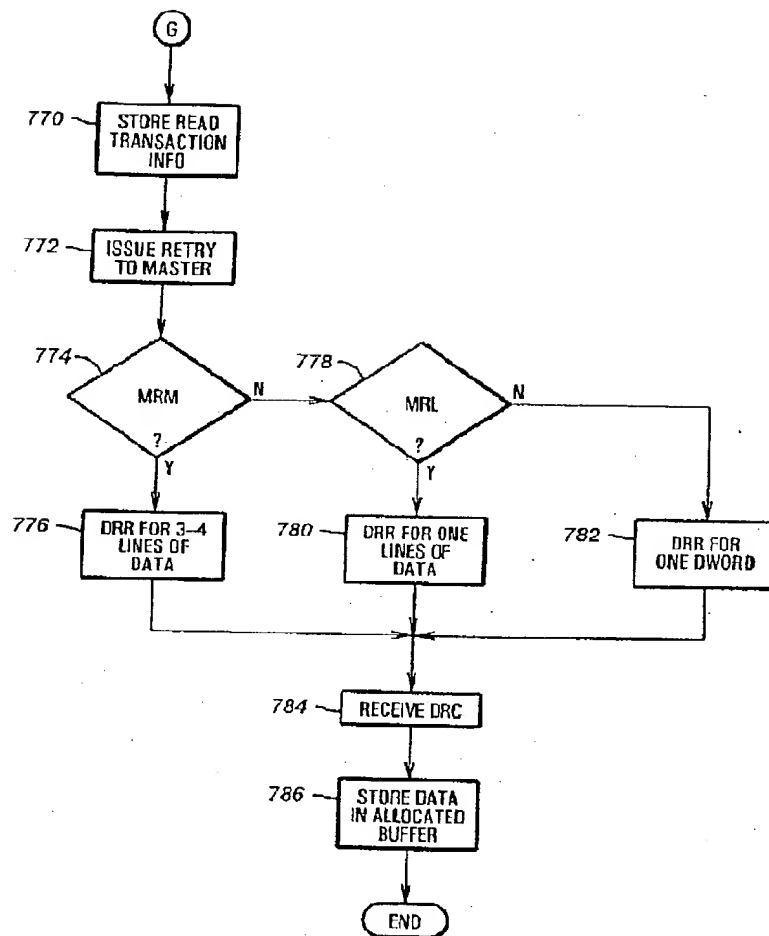
【図14】



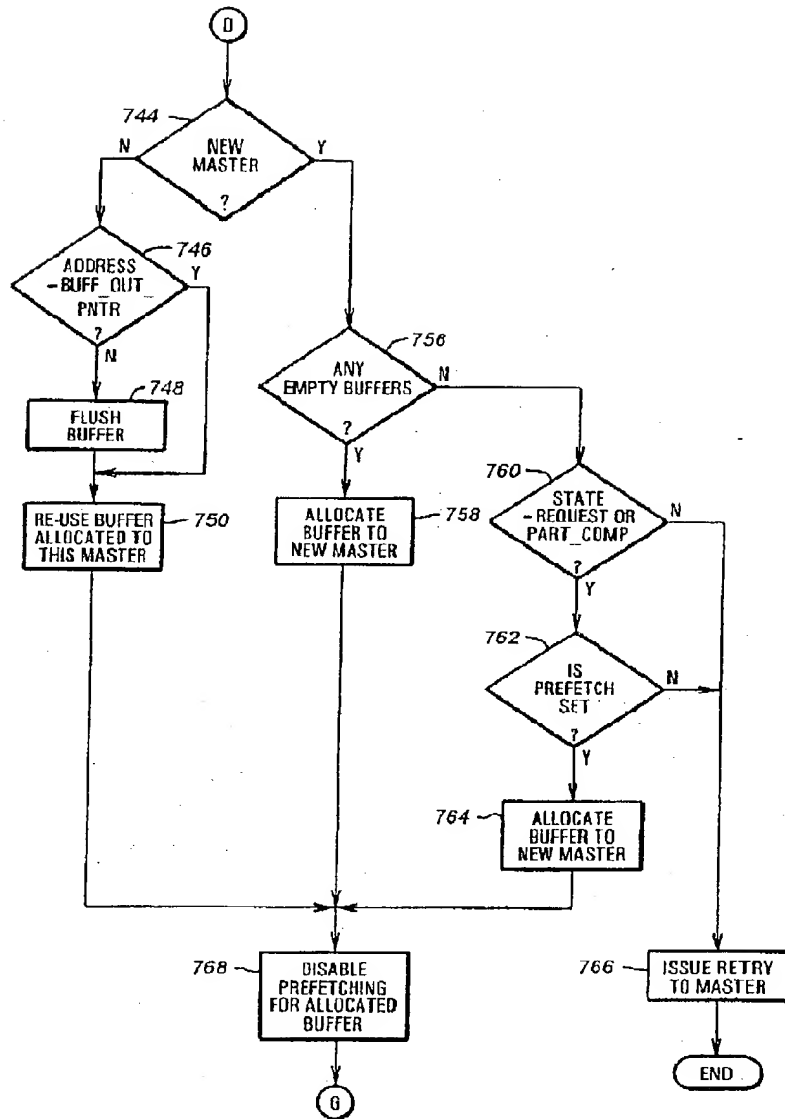
【図15】



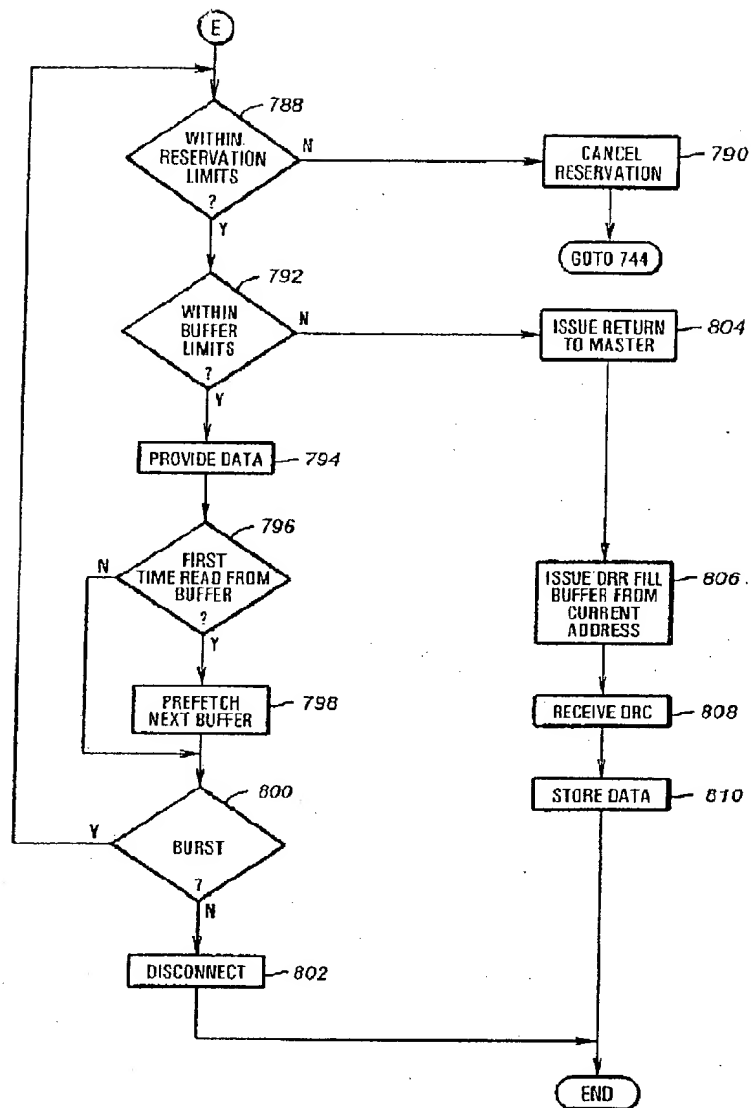
【図16】



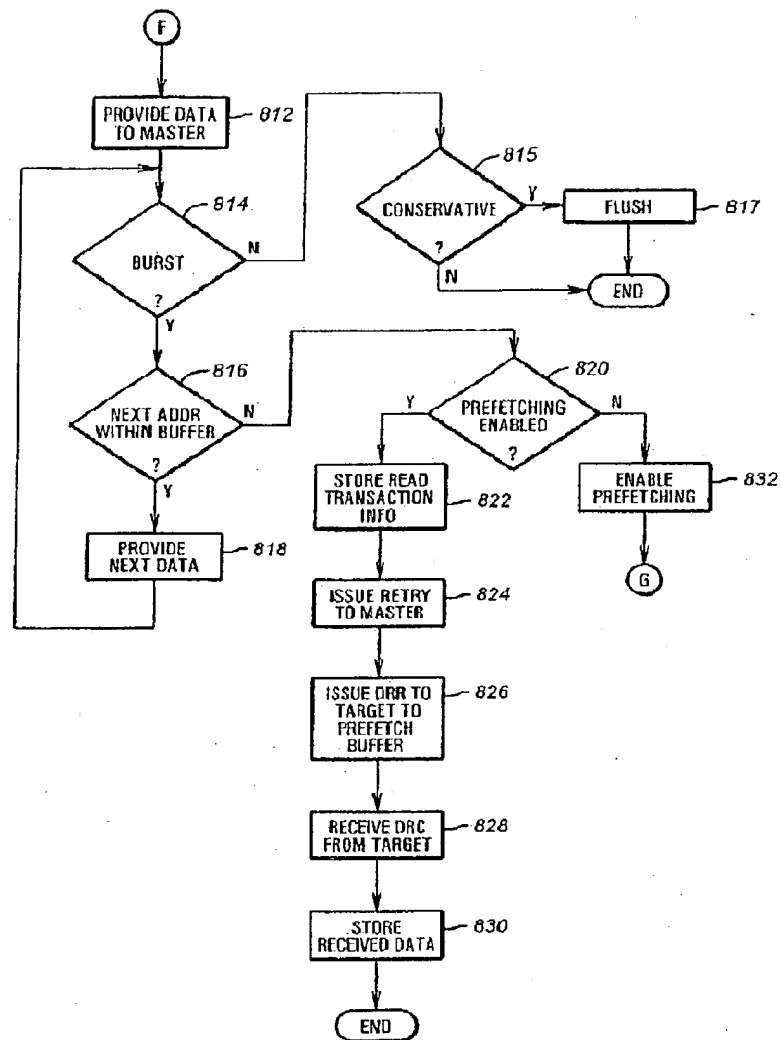
【図1.7】



【図18】



【図19】



1. Abstract

A method for transferring data through a bus bridge. The bus bridge includes a number of data buffers for storing data, prefetching data and write posting data. A device communicating with the bus bridge may reserve a buffer by one of two reservation mechanism. The reservation mechanism provides the bus bridge with the address and byte count. The reservation may also be forwarded to any upstream bus bridges. The reserved buffers are prefetched for efficient use of bus access. Data is prefetched and flushed according to alternative algorithms if a buffer is not reserved.

2. Representative Drawing

Figure 2